

UNITED STATES DISTRICT COURT
WESTERN DISTRICT OF WASHINGTON AT SEATTLE

THROUGHPUTER, INC.,

Plaintiff,

v.

MICROSOFT CORPORATION,

Defendant.

COMPLAINT FOR PATENT
INFRINGEMENT

JURY DEMANDED

ThroughPuter, Inc. (“ThroughPuter”) hereby alleges for its Complaint (“Complaint”) against Microsoft Corp. (“Microsoft”) as follows.

INTRODUCTION

1. A Field Programmable Gate Array (FPGA) is a specific type of microprocessor that can be reconfigured based on the tasks to be performed by it. In some situations, this reconfiguration takes place in between the performance of tasks by the FPGA. For example, the FPGAs can be configured and re-configured to provide hardware accelerators for data processing functions such as encoding/decoding, encryption/decryption or compression/decompression and then reconfigured to perform speech translation in addition to, *e.g.*, encryption and compression. Starting in 2015, Microsoft started putting FPGA processors in each server for Microsoft Azure

1 (“Azure”), which is the world’s largest cloud computing platform-as-a-service (PaaS). Using these
2 FPGA processors for the underlying dynamic parallel execution environment or architecture,
3 Microsoft claims a 150-200 fold improvement in data throughput and a 50 fold improvement in
4 energy efficiency. Exhibit 15 at 91, Exhibit 29. Latency has also lowered by about a factor of 10.
5 Exhibit 15 at 86, Exhibit 29. The end result is a power savings in Microsoft and increased
6 processing speeds for Microsoft and the users of applications running on the world’s largest cloud
7 computing platform.

8 2. In 2013, Plaintiff ThroughPuter disclosed a reconfigurable and dynamic parallel
9 execution architecture running on FPGA processors in writing to Microsoft. Two years later,
10 Microsoft filed a patent application on the same hardware-based fabric (failing to disclose to the
11 Patent Office any information about ThroughPuter’s technology or earlier patent filings).
12 Microsoft was ultimately awarded several patents on this subject matter. For example, in October
13 2020, Microsoft was awarded U.S. Patent No. 10,819,657, which claims that match almost exactly
14 those of ThroughPuter’s U.S. Patent No. 11,150,948, which is the basis for Count 1 of this
15 complaint. In other words, Microsoft was awarded a patent on the same hardware-based fabric
16 claimed in the ThroughPuter patents, wrongly suggesting that Microsoft had invented the
17 architecture underlying Azure where ThroughPuter had already patented that architecture and
18 disclosed it to Microsoft.

19 3. The figure below left is an excerpt from the materials disclosed by ThroughPuter
20 to Microsoft in 2013, disclosing an application load and type adaptive “manycore fabric.” The
21 different colors show different reconfigurable cores (*e.g.*, logic blocks of FPGA processors) being
22 assigned different tasks over time. The figure on the right is a colorized version of figures from
23
24

Microsoft's 2015 patent filing. The Microsoft figures show the same reconfiguration of cores to different tasks over time.

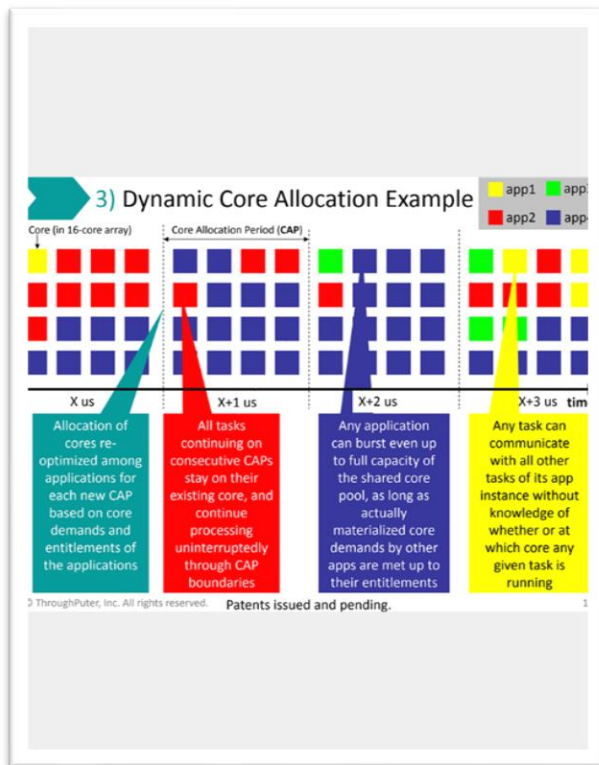


Exhibit 43 at 18.

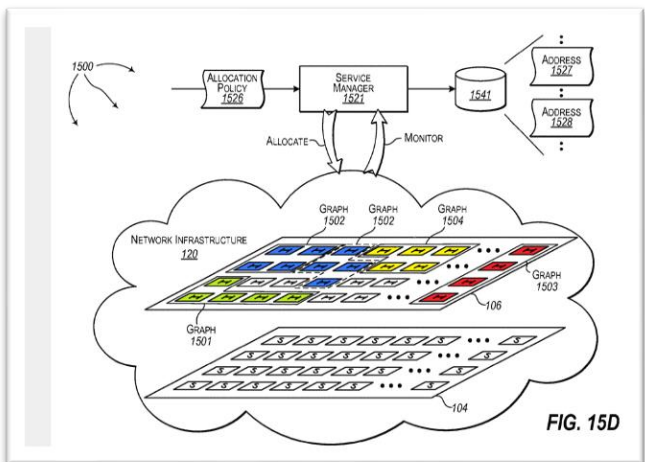
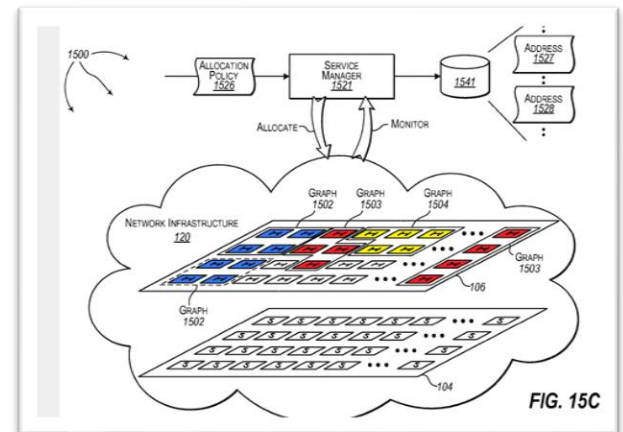


Exhibit 11, FIGs. 15C and 15D (annotated)

4. Microsoft's copying of ThroughPuter's technology was deliberate and wanton. While ThroughPuter has managed to survive as a small company developing and offering for sale other innovative solutions enabled by the throughput improvements which result from ThroughPuter's patented technology, Microsoft's copying and infringement of ThroughPuter's

1 intellectual property has irreparably damaged ThroughPuter's prospects of entering the PaaS
2 market, which had grown from \$3.8B in 2015 to \$37.5B in 2019.¹

3 5. On March 31, 2021, ThroughPuter filed suit against Microsoft in the United States
4 District Court for the Eastern District of Virginia under Civil Action No. 3:21-cv-00216-MHL
5 alleging infringement of related patents (the "Virginia Action"). ThroughPuter served the
6 complaint in the Virginia Action on April 6, 2021. The Eastern District of Virginia transferred the
7 Virginia Action to this Court by Order dated March 23, 2022. This Court assigned the transferred
8 Virginia Action to Judge Rothstein as Civil Action No.2:22-cv-00355-BJR, which remains
9 pending. This complaint asserts two patents that had not yet issued when ThroughPuter filed the
10 Virginia Action.

11 6. Prior to filing the Virginia Action ThroughPuter repeatedly attempted to engage
12 with Microsoft to resolve this matter outside of litigation. Microsoft rejected or ignored
13 ThroughPuter's proposals and requests for business resolution of the matter.

14 7. Microsoft's decision to reject ThroughPuter's proposals and offers to collaborate
15 created a highly inequitable situation by which Microsoft has used its tremendous market power
16 to scale up the world's largest and most successful cloud computing platform. Microsoft did this
17 based on the technology it copied from ThroughPuter: a small start-up that could have only
18 competed based on its effort to protect and patent its innovations. Microsoft vaporized the
19 competitive edge to which ThroughPuter was entitled by its decision to ignore ThroughPuter while
20

21
22
23
24 ¹ <https://www.statista.com/statistics/505248/worldwide-platform-as-a-service-revenue/>

1 generating billions of dollars in revenue per quarter based on its unlawful exploitation of
2 ThroughPuter's technology.

3 8. In a 2016 presentation entitled "Catapult at Ignite Innovation Keynote,"
4 Microsoft's Chief Executive Officer, Mr. Satya Nadella, boasted of the numerous benefits
5 Microsoft had achieved by incorporating FPGAs into its cloud computing offerings: "We now
6 have FPGA support across every compute node of Azure. That means we have the ability, *through*
7 *the magic of the fabric* that we have built, to distribute your machine learning tasks, your deep
8 neural nets, to all of the silicon that is available, so that you can get that performance, that scale."
9 Exhibit 31 at 9 (emphasis added). The "magic of the fabric" did indeed allow Microsoft to scale
10 up an FPGA-based cloud computing solution. But the "magic of the fabric" was not invented by
11 Microsoft, it was copied from ThroughPuter.

12 9. According to Dr. Doug Burger (with whom ThroughPuter corresponded about its
13 technology), the "magic of the fabric" referred to by Microsoft's CEO, use of FPGAs "allows us
14 to do things on a scale that hasn't been done before." According to Dr. Burger: "It gives us the
15 most powerful cloud, the most flexible cloud, and the most intelligent cloud." *Id.* at 30, 34.

16 10. The loss of ThroughPuter's technical competitive advantage through Microsoft's
17 copying based infringement, which Microsoft admits gives it on the order of 100-fold performance
18 and efficiency gain, has devastated ThroughPuter's business, including the ability to raise
19 sufficient startup capital, gain collaborators, partners and initial customers, and to generally enter
20 the market with a differentiated or more cost-efficient solution. ThroughPuter has survived
21 notwithstanding Microsoft's anticompetitive tactics due to its ability to innovate in the fields of
22 application of its core technology.
23
24

11. ThroughPuter now brings this action based on Microsoft's willful infringement of the patents that form the causes of action herein.

NATURE OF THE ACTION

12. This is an action for patent infringement arising under the Patent Laws of the United States, 35 U.S.C. § 1 *et seq.*, including 35 U.S.C. § 271.

13. ThroughPuter brings this action to halt Microsoft's infringement of its rights under the Patent Laws of the United States, 35 U.S.C. § 1 *et seq.*, which arise under the following patents:

- U.S. Patent No. 11,50,948 (the “ ’948 patent), which is attached hereto as Exhibit 1 and
- U.S. Patent No. 11,036,556 (the “ ’556 patent), which is attached hereto as Exhibit 2².

THE PARTIES

14. Plaintiff ThroughPuter, Inc. is a Delaware corporation having a principal place of business at 249 Richmond Road, Williamsburg, VA 23185. Plaintiff owns over 50 issued domestic and foreign patents and pending applications protecting its products, services and technologies. ThroughPuter's President, Mark Sandstrom, is the named inventor on each of such patents and applications.

15. Plaintiff has developed and continues to develop various products and services, including i) Estimator™, a machine learning Application Specific Processor (“ASP”)-as-a-service

² Exhibits 3-8 are intentionally left blank.

1 offering of the ThroughPuter PaaS project, and ii) Grafword™, an artificial intelligence (“AI”)
2 powered, graphical authentication service that is a pilot application of the Estimator™ machine
3 learning microservice.

4 16. Estimator™ provides a streaming machine learning (“ML”) microservice, to
5 support AI applications in unpredictably changing operating environments. Estimator™ allows its
6 prediction models and logic parameters to be adjusted continuously while the microservice is in
7 operation, such that its predictions will stay tuned-in to the prevailing reality of its operating
8 environment, as that may evolve over time or even change abruptly. An International Search
9 Report recently conducted by the International Search Authority under the Patent Cooperation
10 Treaty concluded that this technology is patentable. A beta version of the Estimator™ application
11 programming interface is currently commercially available for 3rd party developer subscription at
12 www.estimatorlab.com.

13 17. Grafword™ provides graphic based high-security password generation and
14 authentication, such that the level of authentication challenge is adjusted according to a level of
15 deviation of a given user’s online session attributes from what is expected for the given username.
16 Grafword™ thus provides both high security as well as, for the authentic users, convenience in
17 online authentication. An International Search Report recently conducted by the International
18 Search Authority under the Patent Cooperation Treaty also concluded that this technology is
19 patentable. A beta version of Grafword™ is used for Estimator™ account creation and login:
20 <https://estimatorlab.com/landing>.

21 18. Both Estimator™ and Grafword™ have the potential to change the space in which
22 they are offered due to the advantages provided by ThroughPuter’s claimed inventions such as
23 increased throughput and latency.

1 19. Microsoft is a corporation organized and existing under the laws of the State of
2 Washington with its principal place of business at One Microsoft Way, Redmond, WA 98052.

3 **JURISDICTION AND VENUE**

4 20. This is an action for patent infringement which arises under the Patent Laws of the
5 United States, 35 U.S.C. § 1 et seq.

6 21. This Court has subject matter jurisdiction at least under 28 U.S.C. §§ 1331 and
7 1338.

8 22. Venue is proper in this District under 28 U.S.C. § 1400(b) because Microsoft has
9 committed acts of infringement and has a regular and established place of business in this District.

10 23. This Court has personal jurisdiction over Defendant Microsoft pursuant to due
11 process and/or Washington's Long Arm Statute because Microsoft has committed and continues
12 to commit acts of patent infringement, including acts giving rise to this action, within the State of
13 Washington and this District, and because Microsoft recruits Washington residents, directly or
14 through an intermediary located in this state, for employment inside or outside this state.

15 24. The Court's exercise of jurisdiction over Microsoft would not offend traditional
16 notions of fair play and substantial justice because Microsoft has established at least the required
17 minimum contacts with the forum.

18 25. Microsoft maintains regular and established places of business throughout
19 Washington and in this District.

20 26. Microsoft has substantial business contacts within this District and has purposefully
21 availed itself of the privileges and benefits of the laws of the State of Washington.

BACKGROUND

27. This case involves ThroughPuter's patented cloud computing, computing acceleration and related technologies, which were developed starting in 2010.

28. As of that time, advancements in computing technologies had generally fallen into two categories. First, in the field conventionally referred to as high performance computing, the main objective has been maximizing the processing speed of a given computationally intensive program running on dedicated hardware. In this field, speed was traditionally achieved by assigning a combination of separate parallel processors to all work on the same program simultaneously. Second, in the field conventionally referred to as utility or cloud computing, the main objective has been to most efficiently share a given pool of computing hardware resources among a large number of client application programs.

29. Thus, in effect, one branch of computing innovation has been seeking to effectively use a large number of parallel processors to accelerate execution of a single application program by parallelizing its processing across a maximum possible number of processors. At the same time, another branch of computing innovation has been seeking to share a single pool of computing capacity among a large number of application programs to optimize utilization of processing capacity. The former efforts pursue maximizing processing speed of a single program. The latter efforts pursue maximizing utilization of processing capacity.

30. As of the time of ThroughPuter's pioneering patent filings starting in 2011, there had not been major synergies between the effort to increase processing speed of a single program on the one hand, and maximizing processing capacity utilization on the other. Indeed, pursuing one of these traditional objectives often happened at the expense of the other, placing the two objectives in tension with each other.

1 31. For instance, while dedicating an entire parallel processor based (super) computer
2 to each individual application would increase processing speed of the individual programs, it
3 would also cause severely sub-optimal computing resource utilization, as much of the capacity
4 would be idle much of the time. On the other hand, while seeking to improve utilization of
5 computing systems by sharing their processing capacity among a number of applications would
6 lead to enhanced resource utilization, it also tended to slow down processing of individual
7 programs. As such, the overall cost-efficiency of computing was not improving as much as
8 improvements toward either of the two traditional objectives would imply: traditionally, increases
9 in processing speed came at the expense of system utilization efficiency, while overall system
10 utilization efficiency maximization came at the expense of individual application processing
11 speed.

12 32. The foregoing tension was exacerbated by the fact that even mainstream application
13 performance requirements were increasingly exceeding the processing throughput achievable from
14 a single CPU core, *e.g.*, due to the practical limits being reached on the CPU clock rates. This
15 created an emerging requirement for intra-application parallel processing (at ever finer grades)
16 even for mainstream programs in order to pursue satisfactory processing speeds, while these
17 programs were to be increasingly hosted on cloud platforms where the processing resources would
18 be shared among programs of multiple clients.

19 33. These internally parallelized and/or pipelined (concurrent) enterprise and web
20 applications would ultimately be largely deployed on dynamically shared cloud computing
21 infrastructure by entities such as Microsoft using the technologies patented, pioneered and
22 promoted by ThroughPuter.
23
24

1 34. Given the foregoing, there existed a need as of 2011 for supporting a large number
2 of concurrent applications on dynamically shared parallel processing resource pools. This then-
3 existing need for a new parallel computing architecture could be met by a system that enabled
4 increasing the speed of executing application programs (including through execution of a given
5 application in parallel across multiple processor cores and/or using hardware accelerators) while
6 at the same time improving the utilization of the available computing resources.

7 35. To address these problems, ThroughPuter developed hardware implemented
8 dynamic resource management functionality including a scheduler, placer, inter-task
9 communications and input/output system for use with multicore processor arrays dynamically
10 shared among multiple concurrent applications, preferably to be deployed on FPGA processors.
11 To that end, in this technology approach, the manycore processor array involves a fabric of
12 reconfigurable cores that can be on-demand programmed to supply the needed mix or match of
13 hardware accelerators. ThroughPuter's technology provided a cloud computing solution that
14 enables accelerated processing speeds across multiple application programs while at the same time
15 optimizing processing resource utilization.

16 36. An exemplary embodiment of ThroughPuter's Dynamic Parallel Execution (DPE)
17 Environment™ (sometimes referred to as "DPEE™" or "DPEE™ technology") is depicted in, for
18 example, ThroughPuter's U.S. Patent No. 9,424,090. Exhibit 9. Referring to the figures and tables
19 of the representative '090 patent reproduced below, this preferred embodiment includes an **array**
20 **of core slots 120**. The processors, or cores, are coordinated and coupled together with a
21 **hardware-based manycore fabric managed by the controller**, which allows the client program
22 tasks to be dynamically placed on processor cores of many-core arrays while communicating with
23 each other directly and securely.

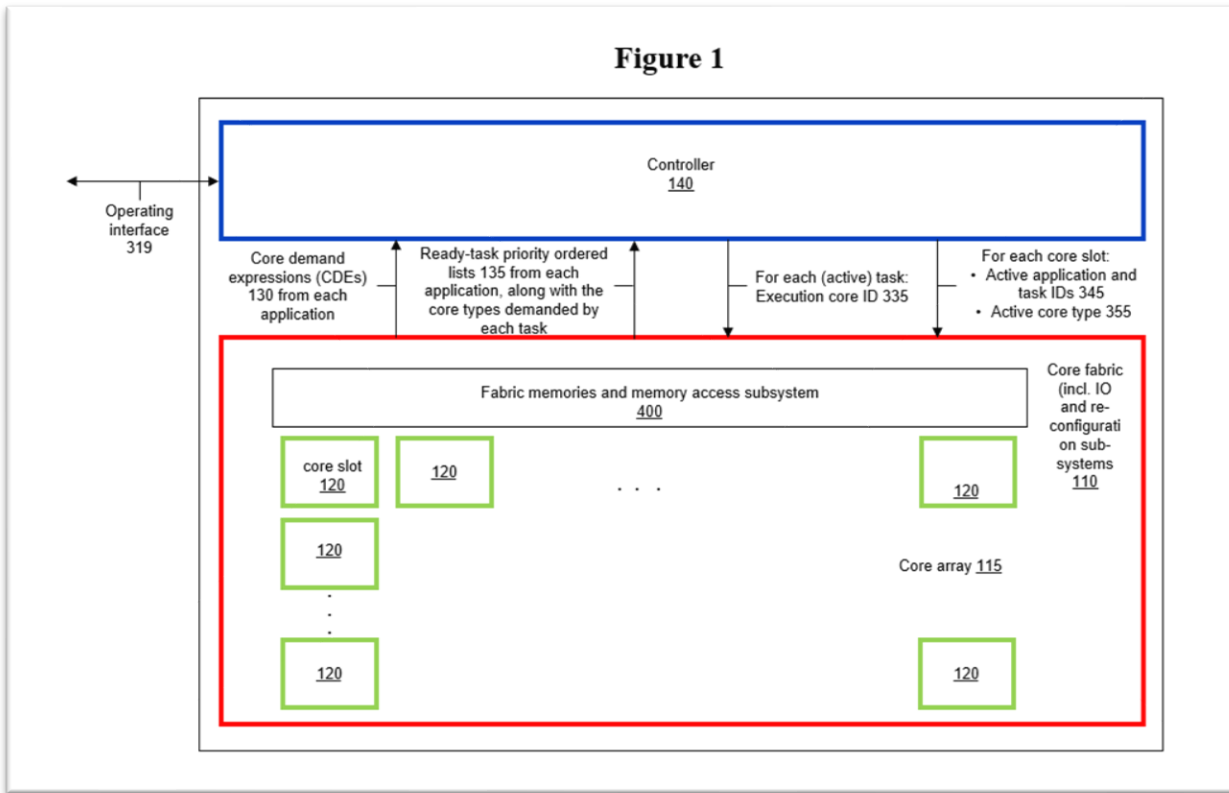


Exhibit 9, FIG. 1 (annotated)

37. ThroughPuter's hardware-based manycore fabric enables processing to be dynamically parallelized and hardware-accelerated, which results in optimized on-time processing throughput across the programs sharing an array of manycore processors. The effect for the client or end user of an accelerated service is increased processing speed and reduced cost base for delivering the application service, such that it becomes economically feasible for cloud service providers to support a range of performance intensive applications even without charge to end-users.

38. In a preferred embodiment, the **controller 140** monitors the processing load for each program sharing **the manycore array 110**, so that periodically re-optimized sets of tasks from the programs can be periodically assigned for execution on the **pool of processor cores 120**. Such periodic reassignment allows for optimal utilization of a pool of processing resources across

a number of potentially competing, concurrent applications. ThroughPuter's exemplary embodiment enables assigning time-variable sets of application tasks for execution on the fabric of reconfigurable cores.

39. A variety of **application tasks 240** are assigned to be run on the **core fabric**, which includes **an array of core slots** within the **core fabric**. Each **application task** may be assigned to run on a **single processor core** within the **core fabric**. The cross-connect, which connects the various components of **the hardware-based manycore fabric**, is used by the **controller** to repeatedly optimize assignment of **tasks** to the appropriate **processor core slots**.

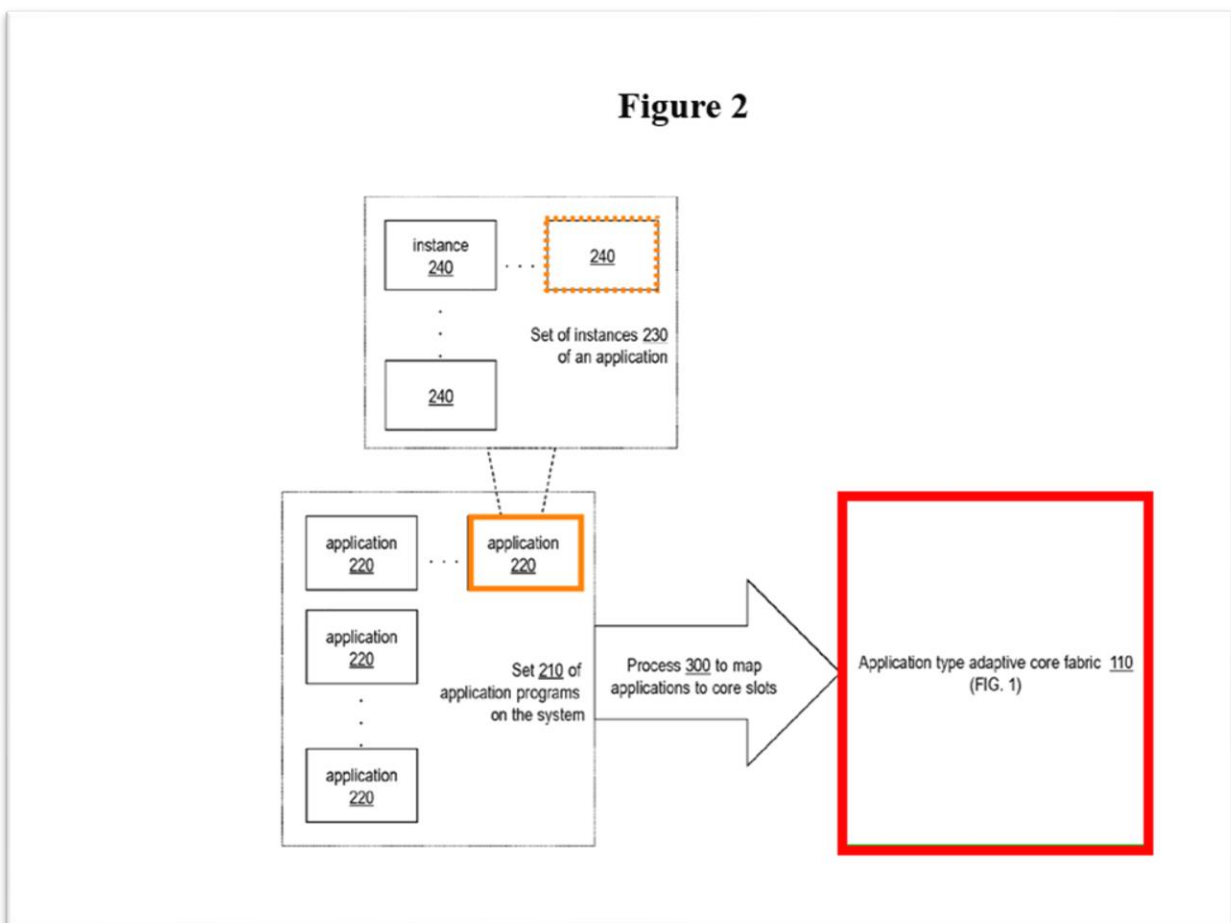


Exhibit 9, FIG. 2 (annotated)

40. The controller manages the mapping of application tasks to processor cores within the core fabric. The result is a system which dynamically and in real-time adapts to varying application processing loads to provide scalable, secure, high-performance and resource-efficient parallel cloud computing.

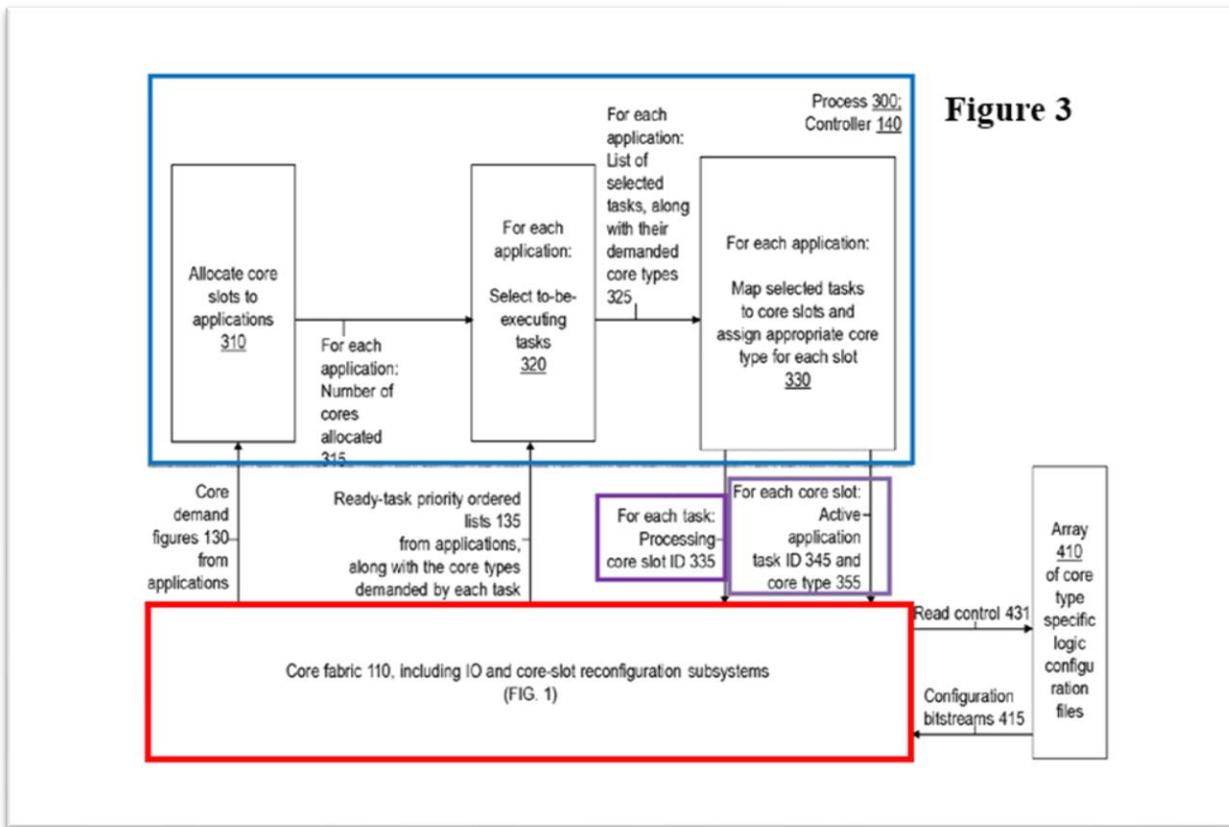


Exhibit 9, FIG. 3 (annotated)

TABLE 5

Core ID index	Application ID	Instance ID (within the application of column to the left)	Core type (e.g., 0 denotes CPU, 1 denotes DSP, 2 denotes GPU, 3 . . . 15 denotes an ASP for a particular function, etc.)
0	P	0	0
1	B	0	0
2	B	8	2
...
14	F	1	5
15	N	1	1

Exhibit 9, Table 5 (annotated)

41. This novel application load and type adaptive manycore fabric permits tasks or applications to be managed at a high degree of granularity with minimized processing overhead while providing various advantages not previously attainable.

42. In recognition of ThroughPuter's innovative achievements, ThroughPuter's Mark Sandstrom was invited to speak at various high performance and cloud computing conferences starting in 2012. GigaOm selected ThroughPuter as one of eleven finalists to present at Launchpad 2012 in San Francisco, CA. That same year, ThroughPuter was invited to present its Dynamic Parallel Execution Environment (DPEE) based PaaS approach at the high performance computing start-up showcase at the Supercomputing 2012 conference ("SC12") in Provo, UT.

43. In February 2013, Mr. Sandstrom reached out to Dr. Doug Burger, Director of Client and Cloud Applications at Microsoft Corporation via email indicating that ThroughPuter

was looking for collaborators. The body of the email, which is attached as Exhibit 19, is reproduced below:

Doug,

I read at http://www.hpcwire.com/hpcwire/2013-02-01/kalray_produces_supercomputer-on-a-chip_for_industrial_applications.html your remarks at HiPEAC, concerning the need to modernize the computing architectures in view of the physical limits and future application requirements. In this context you may find ThroughPuter's cross-layer optimized platform architecture, based on dynamic parallel execution model, quite relevant -- please review <http://www.throughputer.com/platform.html>

ThroughPuter is looking for collaborators for the effort to make the advanced dynamic parallel program execution capabilities available for users (application developers) via PaaS model; possibly Microsoft Azure business unit might be interested in exploring the collaboration opportunities.

Feel welcome to share this call for collaboration among the appropriate parties at Microsoft, and naturally please get back to me for any questions etc. further discussions.

44. Dr. Burger responded "Thanks Mark ... I appreciate the note. I'll forward to the right people in Azure." *Id.*

45. In February 2013, the webpage at <http://www.throughputer.com/platform.html> had the following content at the time of the email exchange:

Platform Overview

ThroughPuter does not build on pre-cloud and sequential processing era concepts such as standalone processor cores and manycore processors as collections of them, or inter-core/process communications or operating systems retrofitted for parallel cloud computing. Instead, ThroughPuter [platform](#) architecture is cross-layer optimized for dynamic, high-performance, high-efficiency, secure cloud computing.

In ThroughPuter architecture, parallel processing hardware is not a mere manycore processor. And neither is its fabric network mere wires and switches between the cores. Nor is its operating system based on conventional sequential OS models with parallel processing as something of an afterthought.

ThroughPuter is a parallel program development and execution platform-as-a-service designed for dynamic cloud computing. For instance, ThroughPuter execution environment is an actual dynamic, secure cloud processor. In ThroughPuter execution environment, the hardware operating system, the adaptable fabric of cores, the fabric network and memory architecture as well as the contract management subsystems all are part of the integrated platform, and work seamlessly together. Collectively, the dynamic parallel processing platform achieves maximized processing throughput, per unit cost, across all client programs dynamically and securely sharing a pool of processing resources.

For illustration, please review technical overview of ThroughPuter PaaS solution:

<http://www.throughputer.com/uploads/ThroughPuterPaaSforParallelProcessing.pdf>

Quick Features -> Benefits -> Customer Value sheet:

<http://www.throughputer.com/uploads/DynamicParallelExecution-FeaturesBenefitsValue.pdf>

Hard-core technology paper on dynamic parallel execution environment of ThroughPuter PaaS:

<http://www.throughputer.com/unloads/DynamicParallelExecutionEnvironment.pdf>

46. The third linked PDF entitled *Parallel Program Execution in a Dynamically Shared Adaptive Manycore Processor* is a document that closely follows the detailed description of ThroughPuter's '090 patent. Exhibit 20.

47. The second linked PDF entitled *Dynamic Parallel Execution via PaaS* is a document that provides an overview of the innovations described in ThroughPuter's portfolio of patents and applications including the application leading to the '090 patent and the substantial advantages they provide. Exhibit 21.

48. The first linked PDF entitled *Parallel Computing Development and Hosting Platform as a Service* is a presentation that provides an overview of ThroughPuter's DPPE. Exhibit 22. As shown in one of ThroughPuter's slides, reproduced below, the cores are dynamically re-assigned to instances or tasks of different applications, which are coded by color.

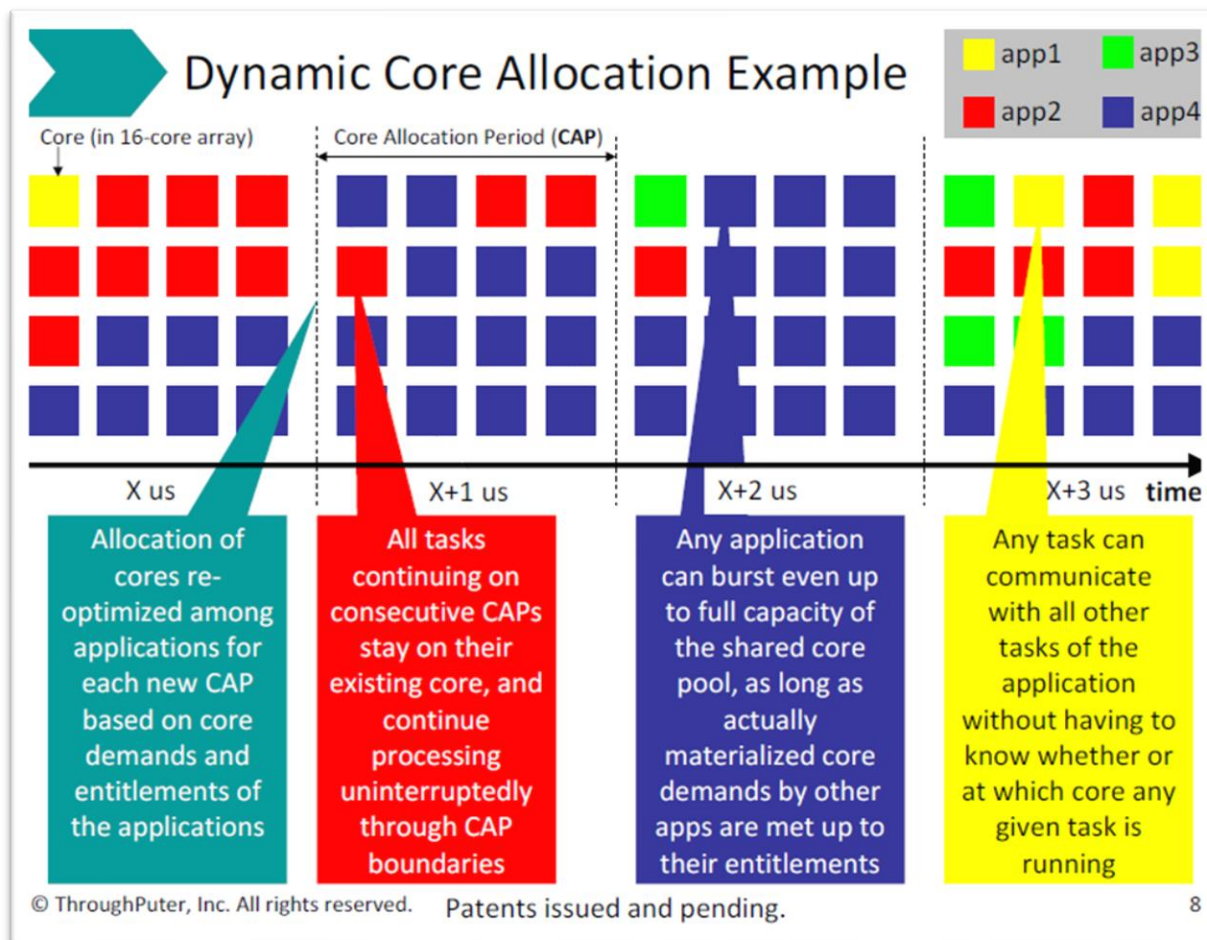


Exhibit 22 at 8

49. Repeatedly, the pool of (sixteen) processing cores are each allocated among the four different application programs. In an illustrative example, these may be application programs used for **speech translation (blue)**, **encryption (red)**, **search ranking (yellow)** and **data compression (green)**. In the first cycle or period, the search ranking (**yellow**) application program is allocated one core while the encryption (**red**) program application and speech translation (**blue**) application program are allocated eight and seven cores, respectively. For the next period, several of the cores are reassigned to the speech translation (**blue**) program application to accommodate an increase in processing demand made by that application. In the third period, a further core is allocated to the speech translation (**blue**) program application and the data compression (**green**)

1 program application is assigned one core, also in response to demand for those programs. In the
2 fourth period, the cores are more evenly allocated across all four applications because the spike in
3 demand for speech translation has abated, or because of demand spikes of the other applications.
4 In this way, the processing throughput and latency across all the applications sharing the given
5 manycore processor array are optimized. The end result for the application programs is faster
6 processing compared to what would be achievable for equal cost base under non-adaptive capacity
7 allocation or without on-demand acceleration.

8 50. ThroughPuter's novel manycore fabric led to industry recognition of ThroughPuter
9 and its technology. For example, in January 2013, ThroughPuter was invited to publish an article
10 in the Cloud Computing Journal, discussing the PaaS based on novel manycore fabric. Exhibit 12.

11 51. In addition, in September 2014, Mr. Sandstrom presented at the FPGAworld
12 conference in Stockholm, Sweden, on the topic of *Hardware Implemented Scheduler, Placer,*
13 *Inter-Task Communications and IO System Functions for Manycore Processors Dynamically*
14 *Shared among Multiple Applications*. Exhibit 43.

15 52. By the time of the 2014 FPGAworld conference, ThroughPuter had already been
16 granted at least a dozen U.S. and United Kingdom patents protecting techniques enabling the
17 advantages of its Dynamic Parallel Execution Environment™ (DPEE).

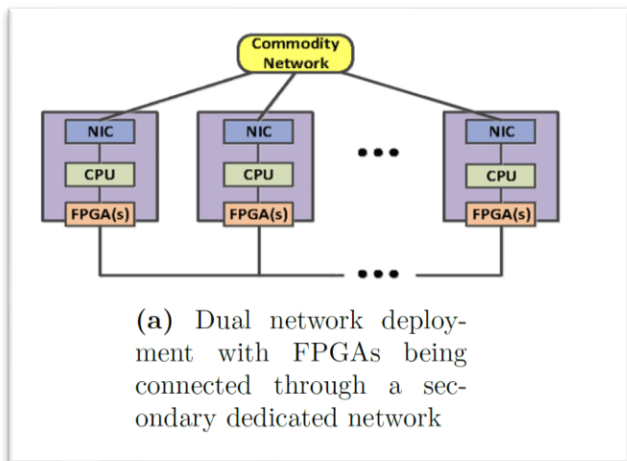
18 53. The following year, ThroughPuter was invited to present at the 2015 HPC Advisory
19 Council Conference in Spain on the topic of executing multiple dynamically parallelized programs
20 on dynamically shared cloud processors. A copy of the presentation is attached hereto as Exhibit
21 24.

MICROSOFT'S INFRINGING CLOUD COMPUTING ARCHITECTURE

54. Microsoft's infringing cloud computing platform is known as Azure which is the technology service platform accused of infringement herein.

55. In April 2014, Microsoft introduced an early version of the Microsoft Azure reconfigurable hardware fabric in the whitepaper *A Reconfigurable Fabric for Accelerating Large-Scale Datacenter Services*, authored by Andrew Putnam, Adrian Caulfield, Eric Chung, Derek Chiou, Paolo Costa, Jeremy Fowers, and seventeen additional Microsoft engineers. Ex.

13. Microsoft initially called this reconfigurable fabric Catapult™ and still uses that term today. *Id.* at 1. This early version of Catapult, depicted in the figure at right, was deployed on a 1,632 server testbed. *Id.*, see also Ex. 40at 28-30.



56. The CPU and FPGA in each node communicate via PCIe. *Id.*, see also Ex. 35passim, Ex. 40at 28-30. *Id.* For external connectivity, dual network interfaces was used. *Id.* The first was Ethernet-based, which enables CPUs to communicate with each other over the commodity data center network. *Id.* The second was a specialized, low latency network implemented using 10Gb SAS cables which connects FPGAs to each other. *Id.* The topology for the second network was a 2D torus. *Id.*

57. The figure below provides details of components within the FPGA. *Id.* Here, custom load logic was implemented within an FPGA fabric allocation called "Role," while the Shell logic was composed of:

- i) PCIe Core: This was used to implement an interface between the host CPU and the corresponding FPGA.
- ii) SerialLite III (SLIII): This was a lightweight protocol used for inter-FPGA communication over the SAS links. Since each FPGA connects to four other FPGAs (2D Torus), there were four SLIII blocks in the Shell.
- iii) DRAM Controllers: These were used to provide on-chip memory access to the Role or host (for DMA over PCIe).
- iv) Inter-FPGA Router: A crossbar connected the four SLIII blocks, PCIe Core and Role.

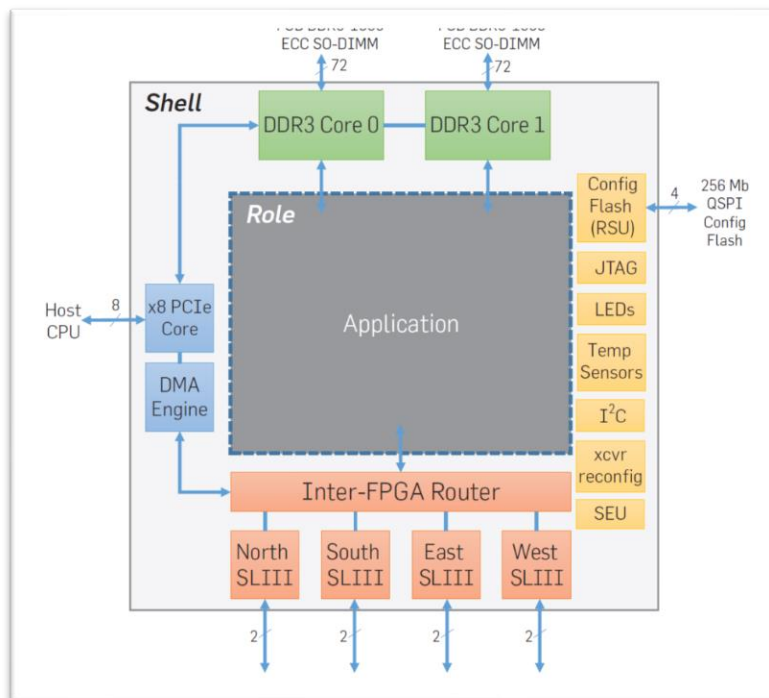


Exhibit 34 at 4; Ex 35 at 5

58. The foregoing version of Catapult will hereafter be referred to as Catapult I.

59. Circa 2015 Microsoft improved the Catapult I architecture to upgrade its initial scope of inter-FPGA connectivity, which permitted only communication between FPGAs in the same rack.

1 60. Whereas Catapult I connected the FPGAs with “complex cabling” that “required
2 awareness of the physical location of machines [i.e., each FPGA]” and provided only for inter-
3 FPGA communication within a single rack (Ex. 14 at 1), the improved version includes a
4 reconfigurable fabric enabling any FPGA in a datacenter to communicate with any other FPGA in
5 the datacenter through super low latency multiplexer-switched connections. *See* Ex. 45(discussing
6 the basics of how a multiplexer serves as a low latency, high speed switch), Ex. 29 at 148, 151.

7 61. This second version of Catapult, called Catapult II hereinafter, has been deployed
8 on Azure servers since the launch of Catapult II circa 2015. Ex. 16 at 6 (indicating deployment in
9 Azure “compute servers beginning in 2015.”); Ex. 10 at 1 (indicating deployment at “hyperscale
10 in Microsoft product data centers worldwide” by October 2016), Ex. 14 at 54.

11 62. At a high level, the improvement of Catapult II involved transforming the network
12 from being a network of interconnected CPUs (some of which might have local, hardwired FGPA
13 acceleration, sometimes called “bolt-on” acceleration) to a network of interconnected FPGAs
14 (each of which can pass packets along to CPUs that are located “behind” the FPGAs). Stated
15 differently, the network was transformed from being a network of CPUs (as was conventional) to
16 being a network of FPGA processors/accelerators (as described in ThroughPuter’s patent filings,
17 which were disclosed to Microsoft years earlier).

18 63. In one of its whitepapers on the topic Microsoft explained that “[a]lthough this
19 change in the network design might seem minor” it is a “major advance” that made a “profound”
20 “impact [on] the types of workloads that can be accelerated and the scalability of the [system].”
21 Ex. 14 at 54. Microsoft explained that this design improvement goes “far beyond just an improved
22 network design” and “can be seen as a fundamental shift in the role of CPUs in the datacenter.”
23
24

1 *Id.* at 58-59. Indeed, the network was changed from a *network of CPUs* with local FPGAs to a
2 *network of FPGAs* with local CPUs. *Id. passim.*

3 64. In a 2016 presentation entitled “Catapult at Ignite Innovation Keynote,”
4 Microsoft’s Chief Executive Officer, Mr. Satya Nadella, boasted of the numerous benefits
5 Microsoft had achieved by incorporating FPGAs into its cloud computing offerings: “We now
6 have FPGA support across every compute node of Azure. That means we have the ability, ***through***
7 ***the magic of the fabric*** that we have built, to distribute your machine learning tasks, your deep
8 neural nets, to all of the silicon that is available, so that you can get that performance, that scale.”
9 Exhibit 31 at 9 (emphasis added). The “magic of the fabric” did indeed allow Microsoft to scale
10 up an FPGA-based cloud computing solution. But the “magic of the fabric” was not invented by
11 Microsoft, it was copied from ThroughPuter.

12 65. According to Dr. Doug Burger (with whom ThroughPuter corresponded about its
13 technology), the “magic of the fabric” referred to by Microsoft’s CEO, use of FPGAs “allows us
14 to do things on a scale that hasn’t been done before.” According to Dr. Burger: “It gives us the
15 most powerful cloud, the most flexible cloud, and the most intelligent cloud.” *Id.* at 30, 34.

16 66. As will be described further below, since the launch of Catapult II the system has
17 continued to evolve in ways that gradually incorporate more and more of the finer details, features
18 and functionality described in ThroughPuter’s patents. Many functions that were performed on
19 the CPUs in the early days of the Catapult II fabric have been gradually shifted to the FPGAs. The
20 Azure fabric (sometimes referred to hereinafter as the “Catapult II network,” the “current Catapult
21 II network” or simply “Catapult II”) strongly resembles or is nearly identical in relevant respects
22 to the hardware-logic (e.g. FPGA) fabric disclosed in ThroughPuter’s patents.
23
24

67. Returning to the “major advance” that is embodied in the Catapult II reconfigurable fabric, central to that advance was the decision to permit every FPGA in a datacenter communicate with any other through a high-speed, low latency, multiplexer switched fabric as described in the ThroughPuter patents. Exs. 10, 15, 36 *passim*. The figure below shows the “shell” of the Catapult II FPGA. Exs. 10, 14. All of the depicted functionality is performed on the FPGA itself, providing high-speed, low latency interconnectivity. *Id.*, *see also* Ex. 46 *passim*, Ex. 16 at 5.

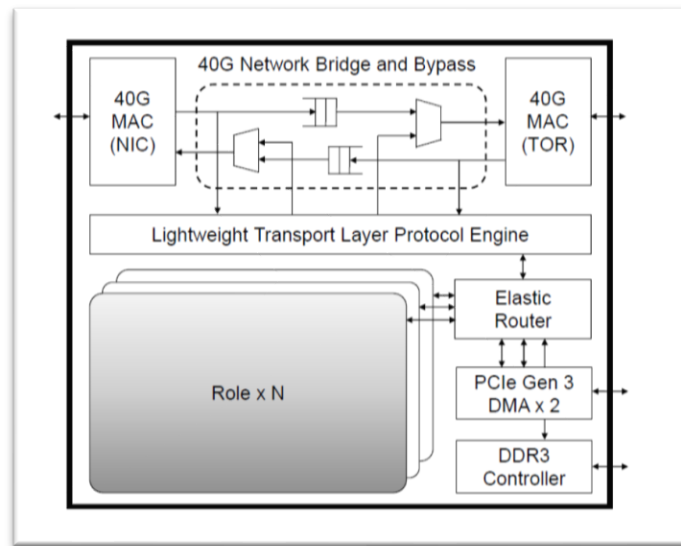


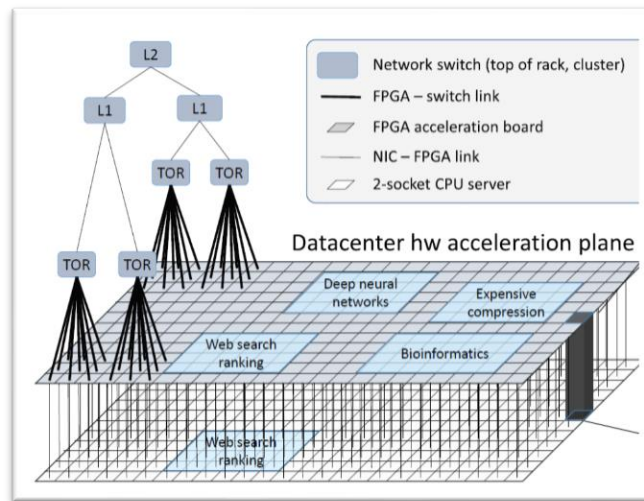
Exhibit 10 at 4

68. Via this shell, the Catapult II inter-FPGA communications and memory accesses occur through the high-speed, low latency multiplexer network inside the FPGA. *Id.* For example, the DDR3 controller comprises DMA multiplexers and demultiplexers. *Id.*, *see also* Ex. 36 at 3. As another example, the bridge between the Lightweight Transport Layer Protocol Engine and the ports for the NIC and TOR includes several functional multiplexer units.

69. In Catapult II, about 23% of each FPGA’s logical circuitry (again, comprising logic units interconnected by high-speed, low latency multiplexer networks) is used to provide the new functionality, including the elastic router, LTL engine, and bridge and the other adapters and

1 controllers of Catapult II. *Compare* Ex. 14 at 4 to Ex. 35 at 5. Collectively, these components
 2 (hereinafter “Catapult II Shell Components”) provide and facilitate FPGA fabric interconnectivity
 3 and memory access. Exs. 10, 14, 36 *passim*.

4 70. Each Azure rack in a datacenter is connected by a two-tier switching network L1/L2
 5 as shown below. Exs. 10, 14, *passim*. The packet communications within each rack and to the top-
 6 of-rack switch (TOR) are facilitated by the FPGA via the high-speed, low-latency multiplexer
 7 controlled connections that serve as the Catapult II Shell Components. *Id*.



16 71. Upon information and belief, since at least 2017, the TOR switch and the L1/L2
 17 switches have included multiplexer-switched network connections to reduce latency. Ex. 38 at
 18 §4.1 (“4.1 Switch design . . . Our circuit switch operates at layer 1, i.e., data traversing the switch
 19 is routed through the PHY block at the ingress and egress ports (Fig. 6). . . The control signals to
 20 these multiplexers are driven by p registers, one per multiplexer. . . Hence, the switch
 21 reconfiguration delay is simply the time it takes to update the registers, which can be done in one
 22 clock cycle.”), Ex. 39 at §4.1 (4.1 Switch Design . . . This is implemented using multiplexers
 23 whose control signal is driven by registers storing the clock counter and timeslot. Therefore, the
 24 switch reconfiguration latency corresponds to the time required to reconfigure the multiplexers,

which is just a few logic gate delays and well below one FPGA clock cycle.”) In 2016, Microsoft published a paper entitled *A Cloud-Scale Acceleration Architecture* authored by, among others, Adrian Caulfield (hereafter “Cloud-Scale Acceleration Architecture”). Exhibit 10.

72. Microsoft Azure includes or has included the functionality as described in Cloud-Scale Acceleration Architecture.

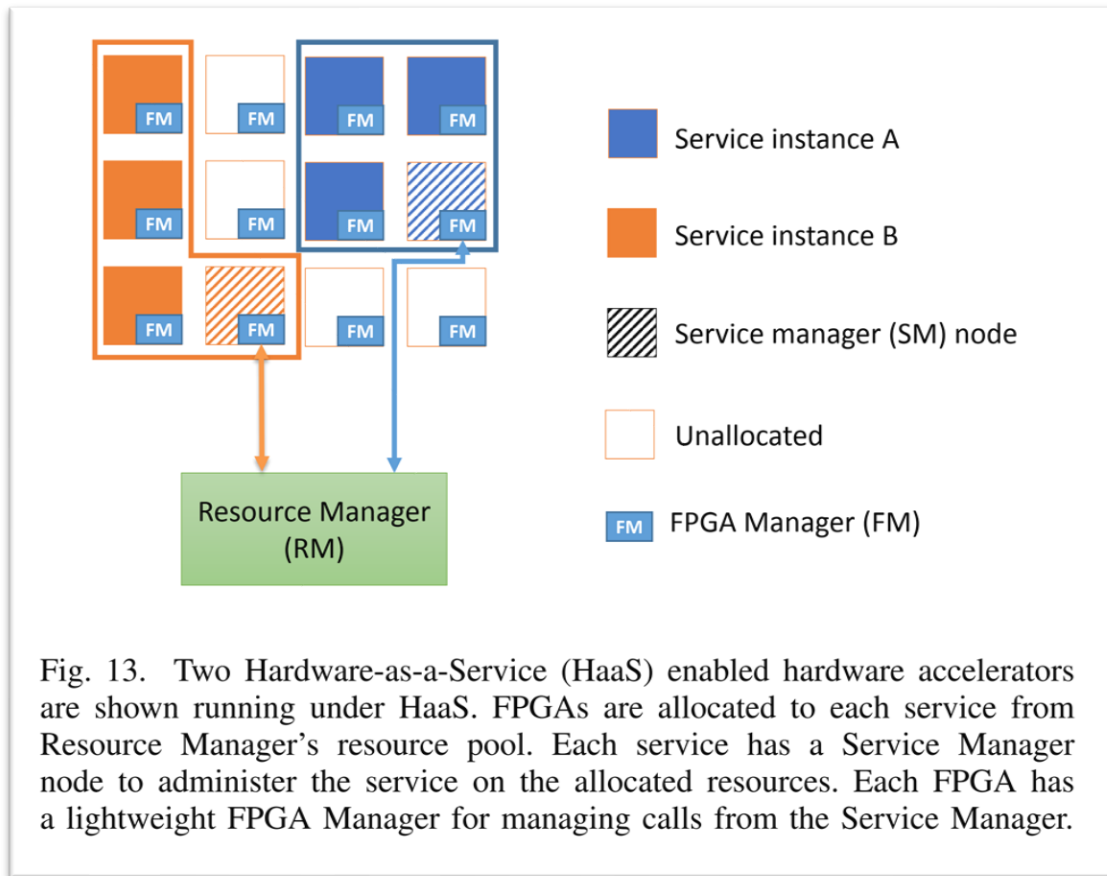


Exhibit 10, Fig. 13

73. Each Service Manager (SM) described in Cloud-Scale Acceleration Architecture is associated with a given, dynamically sized and placed, FPGA processor group, each of which is referred to as a “service instance” in Cloud-Scale Acceleration Architecture.

1 74. Cloud-Scale Acceleration Architecture states that the functionality described
2 therein is used to “accelerate . . . Azure infrastructure workloads,” among other use cases. Exhibit
3 10 at 13.

4 75. Microsoft’s presentation entitled *Inside Microsoft’s FPGA-Based Configurable*
5 *Cloud* (Exhibit 15) demonstrates that Azure embodies or has embodied the functionality described
6 in Cloud-Scale Acceleration Architecture (Exhibit 10) and Microsoft’s U.S. Patent No. 10,270,709
7 (the “ ’709 patent) (Exhibit 11).

8 76. As explained in the presentation, Microsoft Azure is a system comprised of a
9 number of FPGA-accelerated servers, which process tasks from multiple application programs
10 involving a controller that allocates FPGA processors from among a number of FPGA processors
11 and then assigns application tasks to FPGA processors in order to pursue accelerated application
12 processing while seeking efficient processing capacity usage.

13 77. A transcript and corresponding screenshots of Microsoft’s presentation entitled
14 *Inside Microsoft’s FPGA-Based Configurable Cloud* are attached as Exhibit 15. The video is
15 available at https://www.youtube.com/watch?v=v_4Ap1bjwgs.

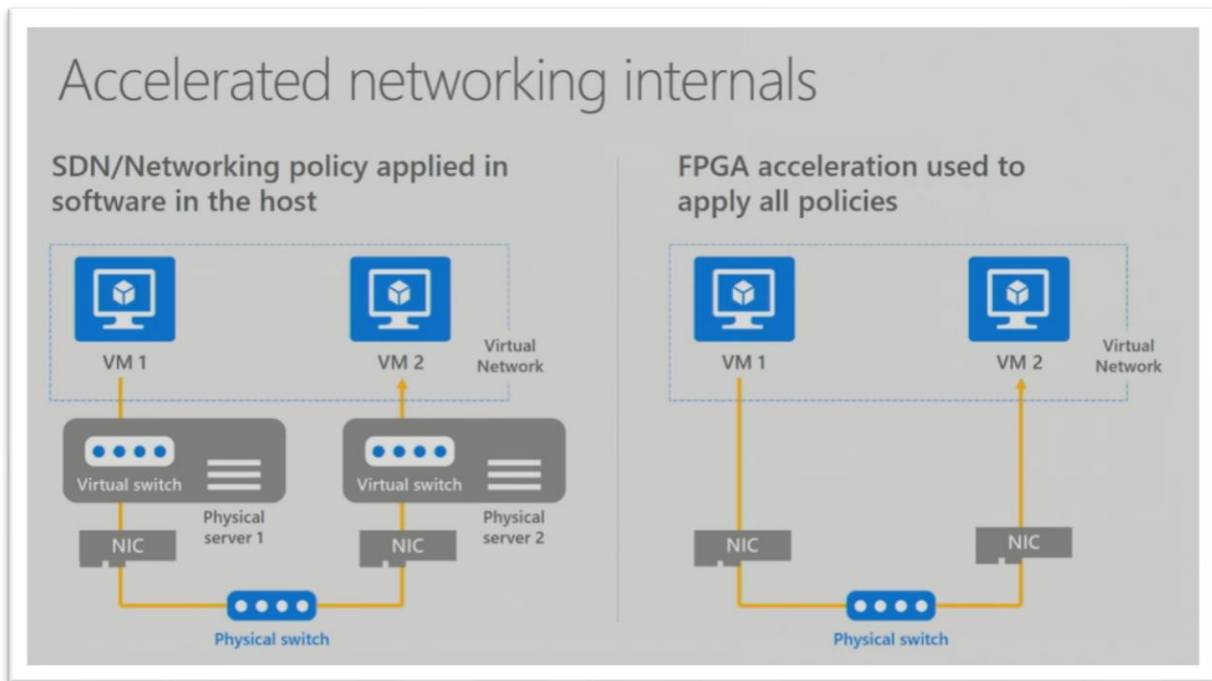


78. Microsoft Azure's RM interfaces with the Azure controller (AC) that applies policies to the data packets transmitted to and from the processors within Azure. Exhibit 15 at 39-42. The RM configures a hardware switch to execute the policies set forth in the hardware access control lists (ACLs), which list grant access rights to a given application program. *Id.*

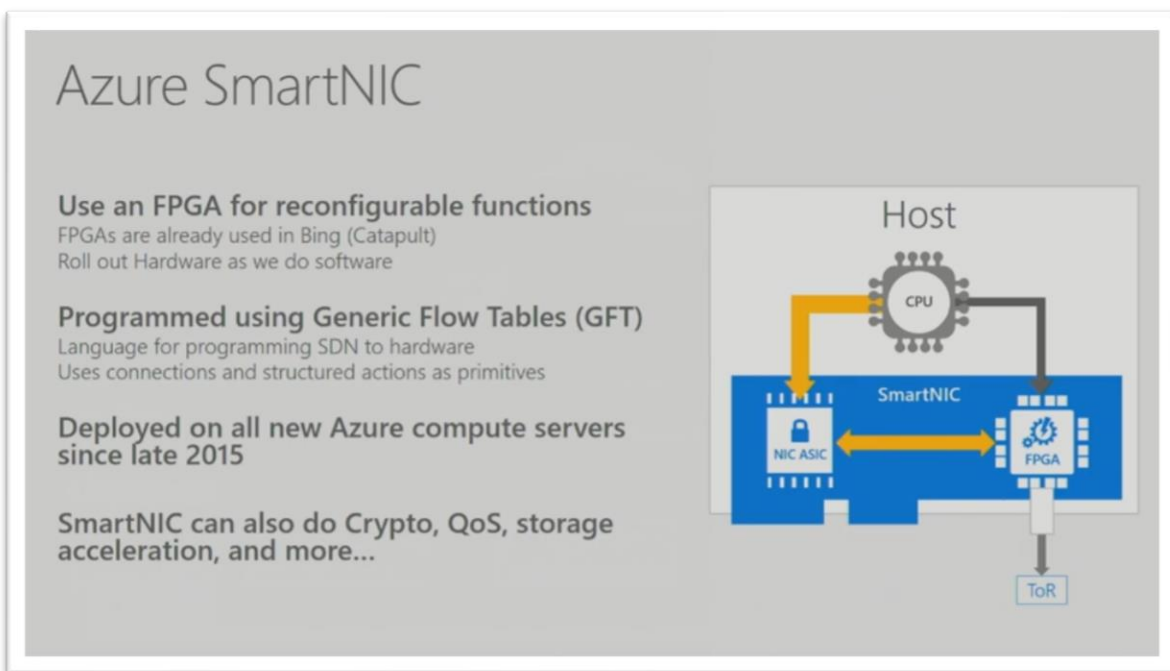
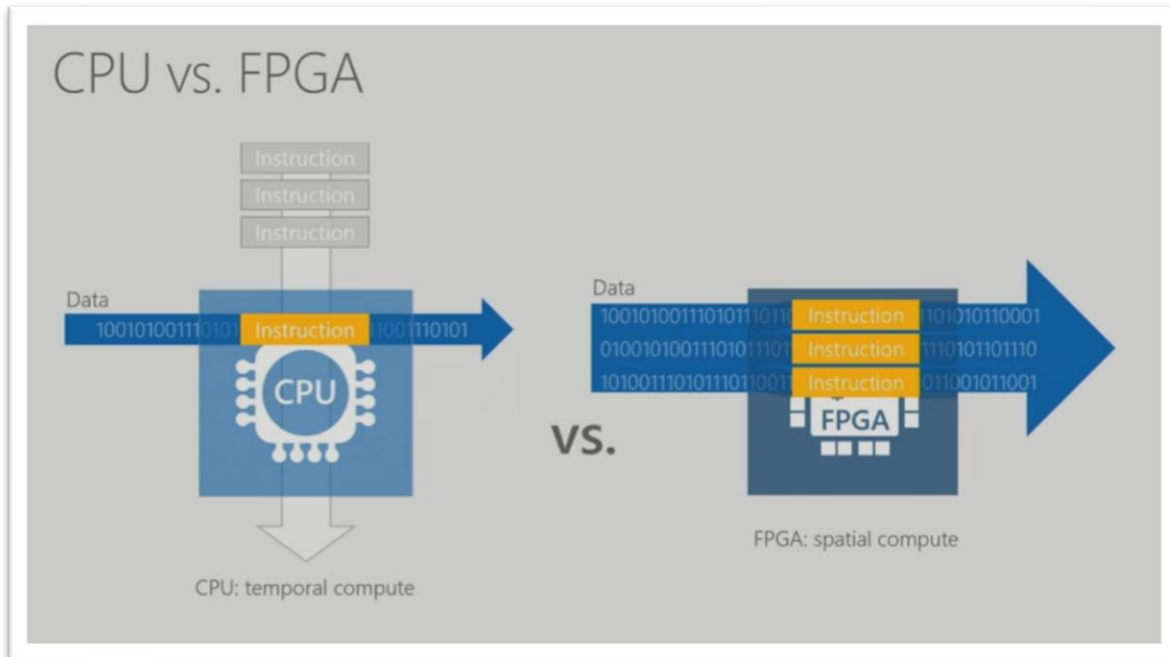
79. Under the direction of the AC, load balancing rules are practiced to pursue optimization of processing capacity by defining how tasks are distributed to FPGA Virtual Machines (*e.g.*, FPGA groupings or graphs in the parlance of Microsoft's '709 patent) (hereafter "FPGA Groupings"). Exhibit 15 at 43-50.

80. The AC instructs the fabric how to route each data packet associated with an application program or service instance task to a given processor based on the type of task to be performed and the available processing capacity on the Azure resource pool. The AC sets up load balancing rules in the load balancing control table (LB NAT) and pushes them out to all servers hosting FPGA Groupings.

81. Under the direction of the AC, Microsoft Azure's physical switch connects directly to the network interface cards (NICs) to route packets (*e.g.*, messages) to application or service instance tasks running on the FPGA Grouping. Exhibit 15 at 54-57. Each FPGA Grouping consists of several FPGA processors operating in a coordinated fashion under the control of the RM in conjunction with an SM, and a FPGA Manager or Node Manager (FM), as explained further below.



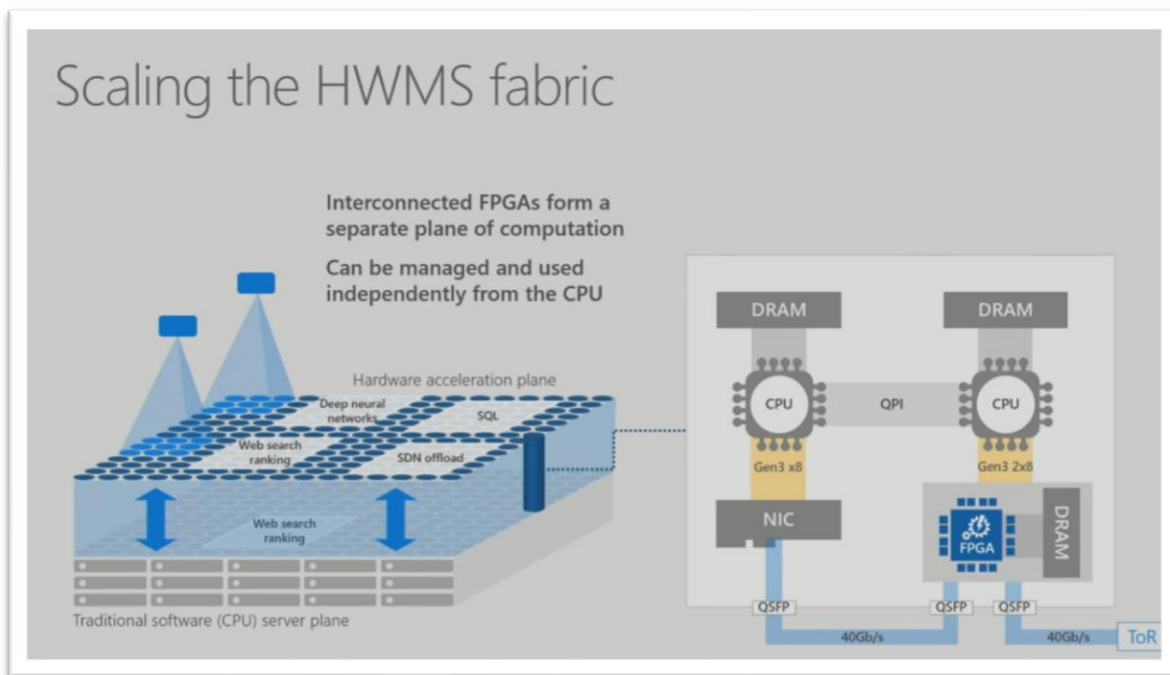
82. The FPGA processors are reconfigured to perform various functions for different application programs. *See, e.g.*, Exhibit 15 at 19-24, 54. As explained further below, the FPGA processors are dynamically reconfigured to be suited for a given task (*e.g.*, deep neural network (DNN) application tasks, SQL database tasks, and so on).



83. As shown in the above figure, an Azure Server includes an FPGA, a NIC, and a CPU. The FPGA is connected with the NIC and CPU, and the Azure Server is connected to other Azure Servers through a network, constituting a fabric of pooled FPGA accelerators. In an

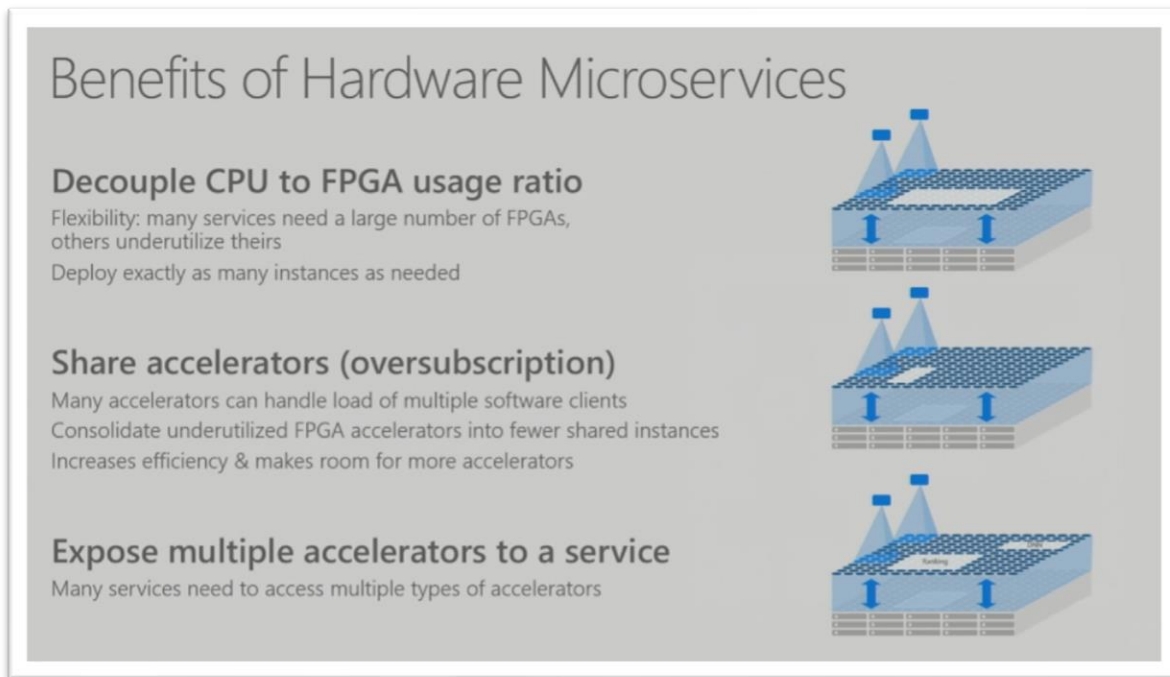
implementation disclosed in Fig. 8 in the '709 patent, an Azure Server includes an FPGA subdivided into separate configurable domains. In this implementation, the separate configurable domains are locally connected to each other through on-chip/on-board hardware. This forms a fabric of accelerator units that may be, at least in part, internal to a subdivided FPGA (hereafter the fabrics referred to in this paragraph are referred to as "Azure Fabric").

84. The FPGA processors are arranged in groups to execute tasks associated with various application programs. Exhibit 15 at 93-98. For example, one FPGA service instance or graph, *e.g.*, an FPGA Grouping, (the blue rectangle with downwardly projecting triangle) executes deep neural network (DNN) application tasks. *Id.* Another executes SQL database tasks, yet another executes network protocol offload tasks, and the last executes web search ranking tasks. Each FPGA Grouping includes several individual FPGA processors (illustrated as blue ovals in the "hardware acceleration plane").



85. In this manner, the AC repeatedly and dynamically rearranges task assignment to

the array of processing units (e.g., FPGA processors) while rearranging communication path connectivity for the array of processing units to optimize the application processing performance as well as the usage of processing capacity on the Azure system. Exhibit 15 at 95-99. The FPGA Groupings (and thus the associated connectivity) changes over time in reaction to the processing demand of the various application programs.



86. Microsoft is now using this FPGA fabric architecture in every new server it deploys in its data centers. Exhibit 28 at 8. Using this architecture, Microsoft obtained 40-100 fold performance improvements in Microsoft Bing's machine learning algorithms. *Id.* Microsoft has also announced the availability of Brainwave, an FPGA-based system for ultra-low latency deep learning for Azure. *Id.*

87. Using this architecture, Microsoft claims up to a 150-200 fold improvement in data processing throughput and up to a 50 fold improvement in energy efficiency. Exhibit 15 at 91, Exhibit 29. Latency has also lowered by about a factor of 10. Exhibit 15 at 86, Exhibit 29.

1 88. In Azure, the RM works in conjunction with both (a) the FM on each host
2 processing unit, which keeps track of resource allocations on the FPGA processors, and (b) the
3 SM, which assigns computing tasks to the FPGA processors. Exhibit 15 at 103-06.

4 89. In the implementation illustrated below, an SM in conjunction with the RM assigns
5 one grouping of computing tasks of a ranking service application program (*e.g.*, Bing) to the FPGA
6 processors that have been allocated. Another grouping of FPGA processors is performing tasks
7 for the Azure Data Link Analytics (ADLA) application, with those tasks having been assigned by
8 another SM in conjunction with the RM, and so forth. In this way, the RM, SM and FMs ensure
9 that each FPGA processor is receiving data from the appropriate queue and executing a task
10 corresponding to the proper application. The RM, SM and FMs also ensure that the output of each
11 FPGA accelerator is communicated to the proper output queue. The load balancing table is utilized
12 to allocate application loads among FPGA processors and FPGAs among loads in order to pursue
13 optimal usage of available processing capacity. In the illustrated example, the Ranking Service is
14 expressing higher processing demand and/or has higher prioritization resulting in it getting
15 allocated four FPGA cores and/or FPGA Groupings, compared with *e.g.*, ADLA that is assigned
16 two FPGA accelerators.

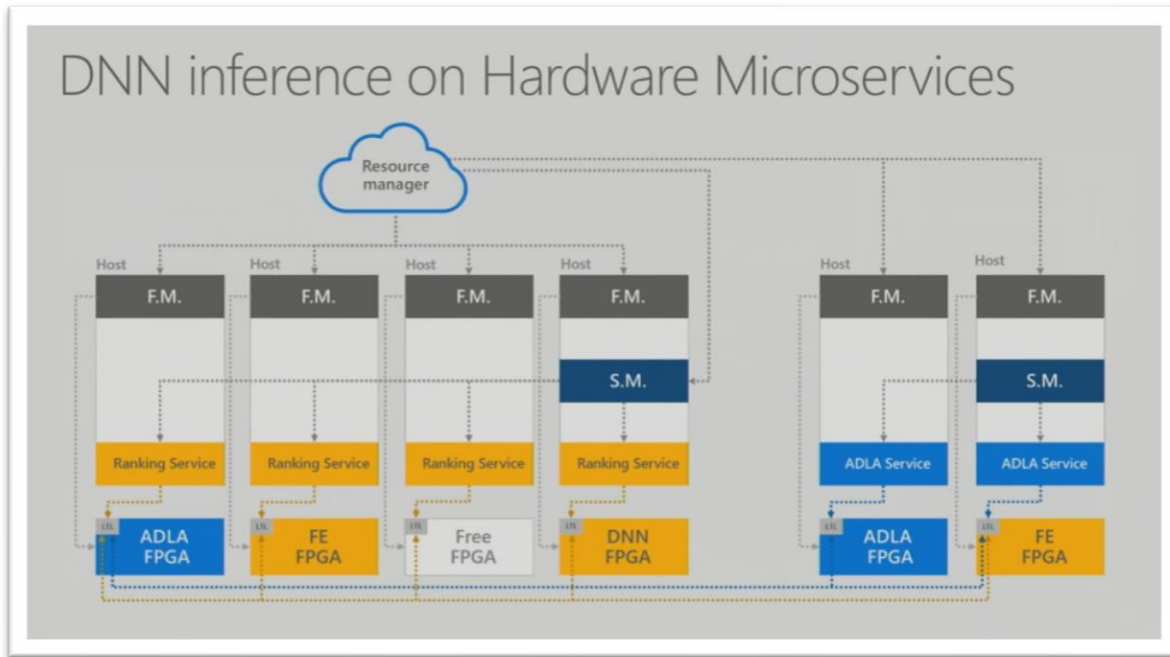


Exhibit 15 at 104

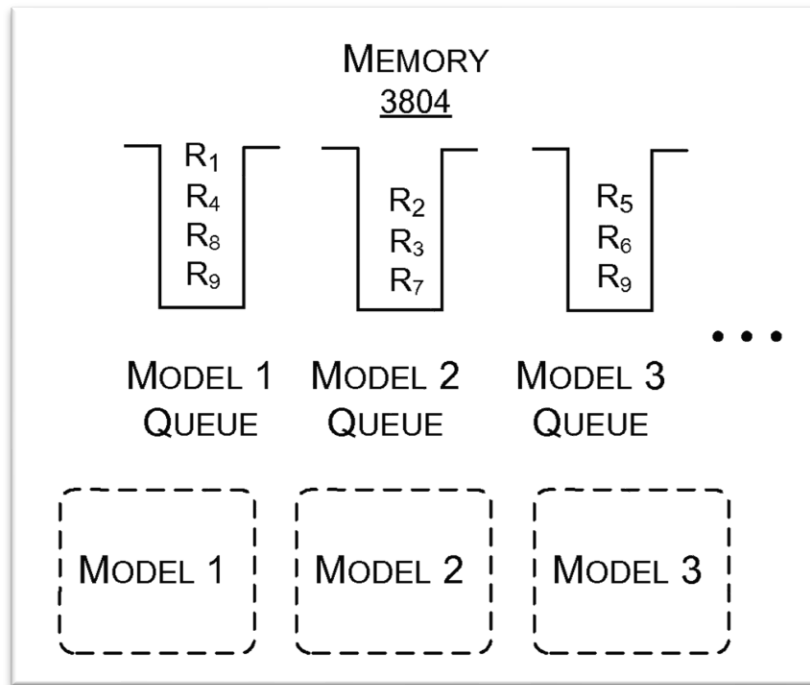
90. Azure's queuing of data as it moves between memories and the FPGA processors is or has been accordant to descriptions in Microsoft's United States Patent No. 10,296,392 ("the '392 patent"). Exhibit 47.

91. On information and belief, Microsoft Azure embodies or has embodied the functionality substantially as described in the '392 patent.

92. In connection with Fig. 38 (excerpted below), the '392 patent explains that different application programs such as a French language search engine, an English language search engine, and a German language search engine may have dedicated input buffers Model Queue 1, Model Queue 2 and Model Queue 3, respectively. Exhibit 47 at 34:62-35:34; see also Fig. 38. Data is fed from the queues to the FPGA processors based on policies such as queue fullness. *Id.* Generally, the queue manager (in coordination with at least the SM) will prioritize the processing of queries in the queue with the most queries. As an example, shown below, Model 1 Queue has

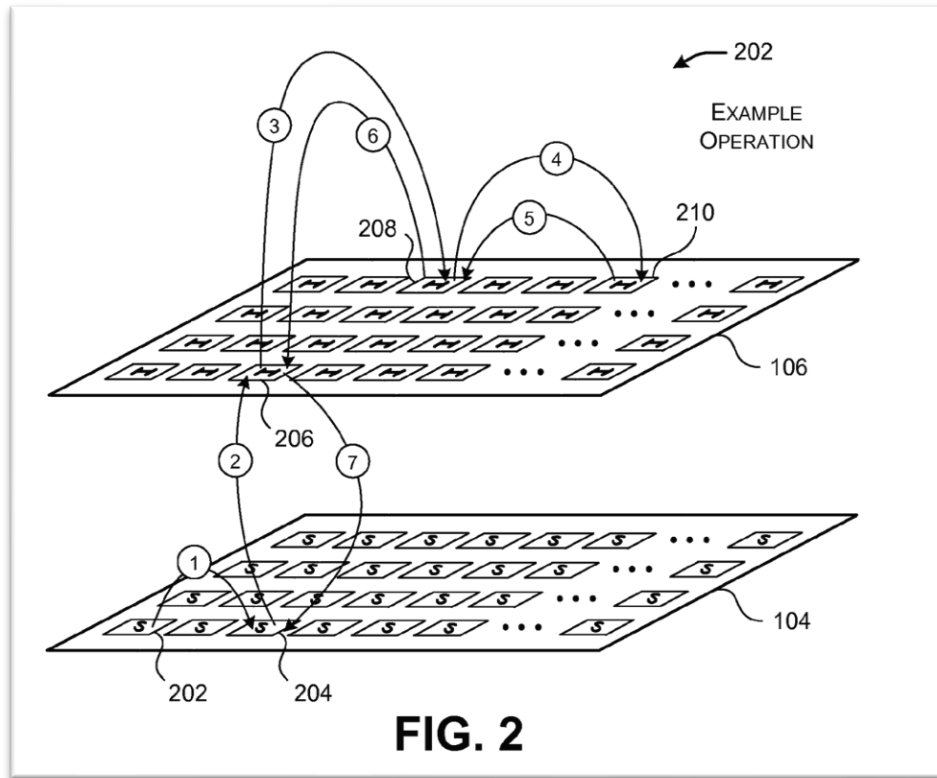
1 queued four queries (R₁, R₄, R₈ and R₉) where Model 2 and Model 3 Queues have queued three
 2 queries. As such, Model 1 Queue would be prioritized above the other two.

3 93. On information and belief, Azure uses such queues to manage the transmission of
 4 data to, from and/or within the FPGA processors.



15 94. The distribution of tasks among the FPGA processors in Azure by the SM, RM,
 16 and FMs is also substantially described in the '392 patent.

17 95. In operations 1-6 illustrated in Fig. 2 of the '392 patent, each hardware acceleration
 18 component 208/210 (which may be an FPGA or a subdivision thereof) performs a task on the data
 19 and forwards the result to the next destination under the management of the head component of
 20 the FPGA Grouping (also referred to as the SM). The '392 patent further explains in connection
 21 with Fig. 36 that each acceleration component performs its respective function, such as
 22 mathematically combining feature values or compressing the data reflecting the feature values
 23 computed thus far.



96. The scheduling and placement of tasks by Azure through the RM, SMs, and FMs is done taking into consideration a given task's readiness for execution. Exhibit 47 at 35:20-25; Exhibit 17 at pp. 287-88.

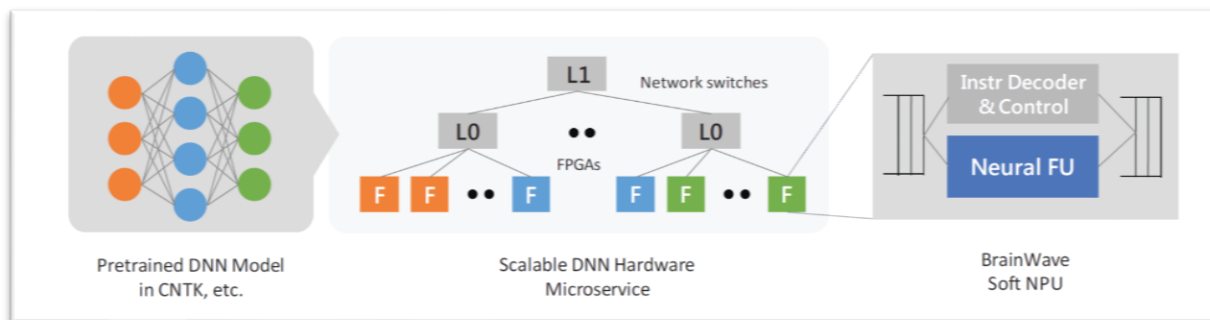
97. On information and belief, Azure embodies or has embodied the technology described in Exhibit 17.

98. Microsoft Azure also supports assured resource allocations for a client's FPGA Grouping for an additional fee. *See* Exhibit 48 at p. 3. Azure allocates FPGA processors based in part on the number of FPGA Groupings subscribed to by the client for its applications.

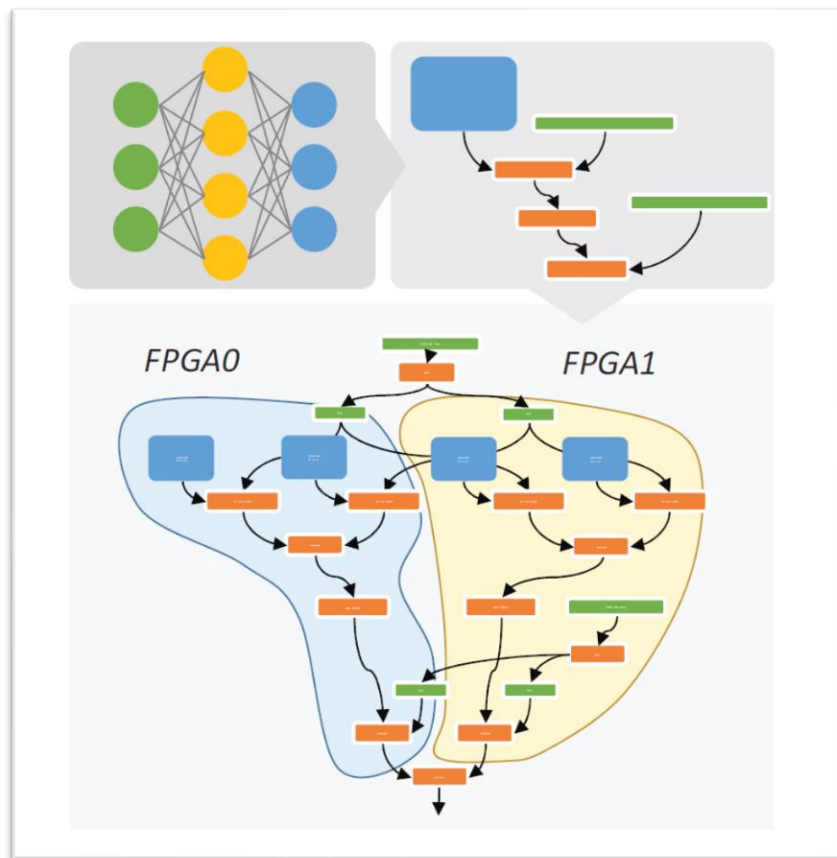
99. In late 2017, Microsoft began moving its deep neural networks (DNNs), which provide the backbone for their real time artificial intelligence (AI) applications, from CPU server networks to Catapult II FPGA server networks. Exs. 18, 37*passim*, see also <https://www.microsoft.com/en-us/research/blog/microsoft-unveils-project-brainwave/>. These

FPGA-enabled DNNs were used internally by Microsoft to support Bing and Azure in 2017 and became available to customers in 2018. Ex. 18 at 17. Microsoft uses the code-name Project Brainwave to refer to its platform that accelerates DNN inferencing, which runs onto of the Catapult servers in Azure.

100. Microsoft's DNN models deployed on FPGA servers (i.e., Brainwave deployed on the Catapult II FPGA network utilize the multi-stage architecture described in various ThroughPuter patent applications, which were published and/or issued as patents more than three years earlier. See, e.g., Ex. 8. As illustrated and explained in one of Microsoft's whitepapers on the topic (Ex. 18), the DNNs are implemented as a hardware microservice having multiple stages. The first stage (illustrated in orange) is implemented on a flexible number of FPGA processing cores and passes its intermediate processing result to the second stage (illustrated in blue), which in turn passes its result on to the last stage (illustrated in green), which outputs the DNN hardware microservice result to the calling application.

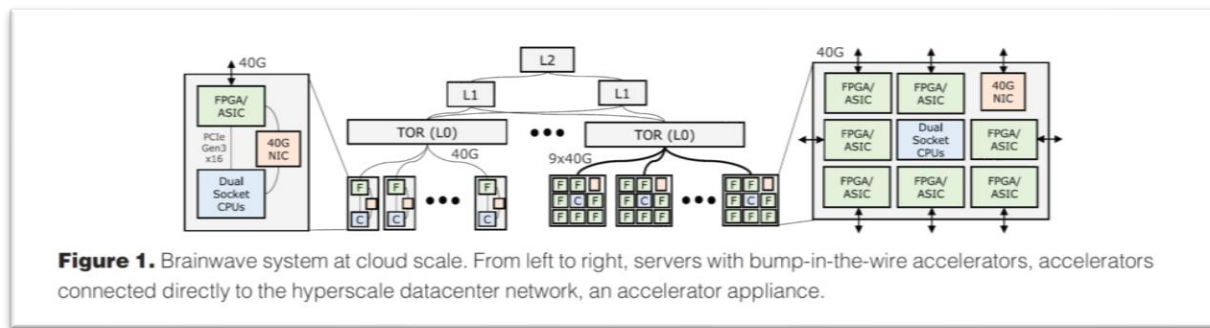


101. The FPGA cores of the different FPGAs that are performing the same task constitute a stage. *Id.* In the illustration below (an excerpt of Ex. 18 Fig. 3), the tasks by the DNN model are separated into three types of tasks, illustrated in blue, orange and green. *Id.* Each of tasks is deployed on a FPGA processing core labelled FPGA0 and FPGA1. *Id.* The blue cores operate collectively to perform a first task in a first stage and pass the output to another a stage comprised of orange cores. *Id.* The green cores operate collectively to perform a second task in a second stage and pass the output to one of a set of stages comprised of orange cores. *Id.* The orange cores perform a third function and serve as the fourth, fifth and sixth stages, although those stages could be considered a single stage for certain purposes and in certain implementations. *Id.*



102. In the Azure-implemented Project Brainwave, an arbitration controller or “network” manages data movement among the memory components: pipeline register files (MRF and VRFs), DRAM, and network I/O queues. Ex. 49 at 6. The arbitration network is deployed on the hardware logic of the FPGA fabric. *Id. passim*. This reduces observed latency by up to 90X. *Id.*

103. In 2018 or 2019, Microsoft began incorporating multiple FPGAs onto each server blade. Ex. 37 at Fig. 1, Ex. 38 at Fig. 1. In this iteration of Catapult II, on each server blade a single CPU is placed “behind” and supports several FPGAs and/or ASICs. This is yet another step in the evolution of Azure away from conventional CPUs and toward more widespread use of the FPGA network to provide an ever-increasing fraction of Azure’s functions.



104. In August 2021, several months after the filing of the First Action, one of the lead Azure engineers, Derek Chiou, published a paper describing a Persistent Recurring Neural Network (RNN) deployed on the FPGA-enabled Azure network. Ex. 23 at 6, 23. Dr. Chiou is an author on most of the whitepapers concerning Azure, including Exs. 10-13, 14, 16, 18, and 36. On information and belief, Microsoft is implementing the technology described in Exhibit 23 on at least some Azure servers.

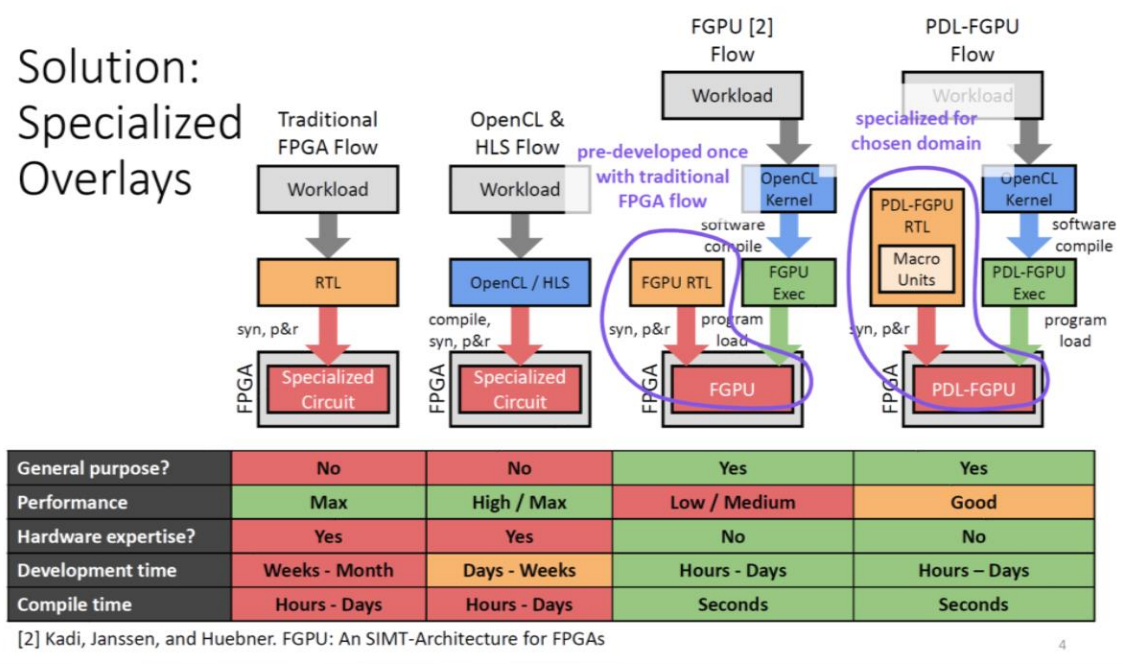
105. RNNs are used for tasks such as speech recognition, text processing and translation. *Id.* Persistent RNNs can be considered a specialized type of DNN in which the inter-stage

1 connections are not static but rather change over time and operations may flow upstream, or
2 backwards through the stages. The latter means that information persists because it can travel
3 backwards through the stages and affect the AI model present at an ostensibly upstream node that
4 is also processing data that was received later in time.

5 106. In its new implementation of the persistent RNN, Microsoft is using essentially a
6 modified shell for the FGPA³. Microsoft's new RNN uses a graphical processing unit (GPU)
7 instantiated on the FPGA. For this reason, it is referred to as an *FPGA general purpose Graphical*
8 *Processing Unit*, or FGPU. Ex. 23 at 4, Ex. 44. Because it is used to perform persistent deep
9 learning (PDL), it is termed a PDL-FGPU. On information and belief, Microsoft currently deploys
10 PDL-FGPUs on some servers to provide Persistent RNNs on Catapult II.

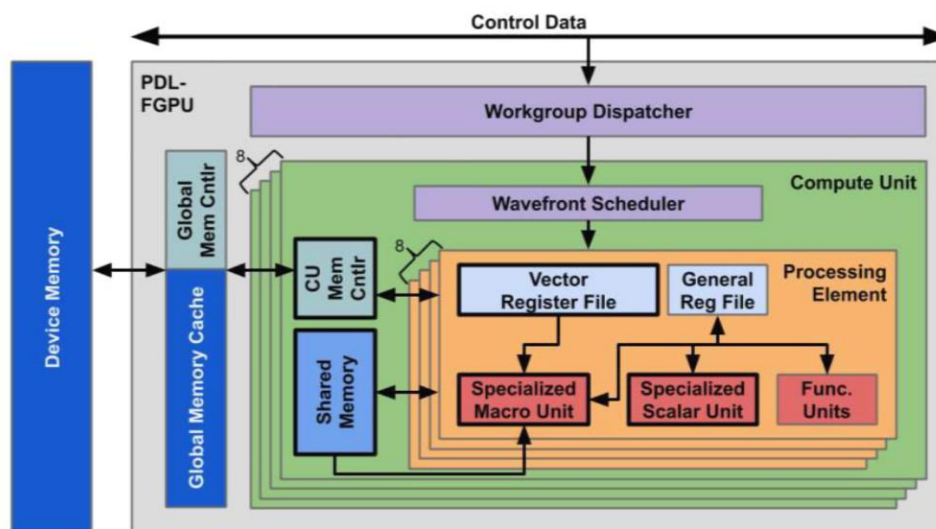
11
12
13
14
15
16
17
18
19
20
21
22
23
24 ³ Recall that shells define the functionality that the FPGA's programmable network of multiplexer-switched ALUs are to perform.

Solution: Specialized Overlays



107. The PDL-FGPU architecture or shell is shown below. As illustrated, the shell controls memory access, including both SRAM on each FPGA and DRAM on each server blade, and communication between each Processing Element (core) in the Compute Unit.

Architecture



1 108. On information and belief, the PDL-FGPU architecture or shell is deployed in
2 conjunction with the Catapult II Shell Components. The PDL-FGPU rides “on top” of the
3 Catapult II Shell Components, thereby allowing the Catapult II Shell-Components to mediate inter-
4 FPGA communication and memory access. On information and belief, Microsoft currently uses
5 this approach on at least some servers.

6 109. Deploying the RNNs on the FPGA-enabled Azure network via PDL-FGPUs
7 provides a one to three order of magnitude improvement in throughput over Microsoft’s internal
8 baseline, which on information and belief is an implementation on a CPU network. See Ex. 23
9 pp. 13-17. On information and belief, Microsoft currently deploys PDL-FGPUs on at least some
10 FPGA-enabled Azure network servers to provide Persistent RNNs.

11 110. To summarize, since the introduction of Catapult II in 2015, Microsoft has
12 transitioned a broad swath of workloads and the associated management responsibility from CPUs
13 over to the low-latency FPGA network fabric of Catapult II, per ThroughPuter’s patents. In its
14 current form, Catapult II functions as stand-alone, almost entirely independent FPGA fabric that
15 handles the majority of the workload in each Azure datacenter. That independent FPGA network
16 is one of the innovations described in ThroughPuter’s patents that has been deliberately copied by
17 Microsoft, including by filing patent applications wrongly claiming these advances as having been
18 invented by Microsoft. In fact, they were disclosed to Microsoft by ThroughPuter years before
19 Microsoft began filing its related patent applications.

20 111. However, Microsoft was right about one thing – this new FPGA fabric architecture
21 is a “major advance” that has a “profound” “impact [on] the types of workloads that can be
22 accelerated and the scalability of the [system].” Ex. 15 at 54. This design improvement goes “far
23
24

1 beyond just an improved network design” and “can be seen as a fundamental shift in the role of
2 CPUs in the datacenter.” *Id.*

3 112. By using the new FPGA-network of Catapult II, as of 2017 Microsoft achieved up
4 to a 150-200 fold improvement in data processing throughput and up to a 50 fold improvement in
5 energy efficiency. Ex. 10 at 34 and exhibits cited therein. Also by that time aggregate latency had
6 been lowered by about a factor of 10. *Id.* On information and belief, in the ensuing four years the
7 evolutions described above (incorporating yet additional functionality described in the
8 ThroughPuter patents) have improved each of these performance metrics by at least another order
9 of magnitude.

10 113. It’s no wonder that Microsoft’s CEO described the new FPGA reconfigurable fabric
11 as “magic.” *Id.* at 4 and exhibits cited therein.

12 **MICROSOFT’S KNOWLEDGE OF**
13 **THROUGHPUTER AND ITS PATENTS**

14 114. Starting in 2011, ThroughPuter has invested in development and patent protection
15 of its intellectual property consistently, in order to protect the competitive advantage to which it is
16 entitled under the law as the inventor of the technology at issue in this case. Recognizing the
17 promise of its technology and the benefits it could provide to Microsoft, ThroughPuter made
18 consistent overtures to Microsoft in an effort to collaborate.

19 115. For example, as discussed above, starting in 2013, Mr. Sandstrom corresponded
20 with Microsoft’s Director of Client and Cloud Applications (Dr. Burger: a named inventor on
21 Microsoft’s U.S. Patent No. 11,819,657, which is discussed further in Count I below) and others
22 on Microsoft’s cloud computing team concerning a potential collaboration between ThroughPuter
23 and Microsoft. *See* Exhibit 19.

1 116. At least through that correspondence, and no later than February 2013, Microsoft's
2 cloud computing team was apprised of the technical details of ThroughPuter's Dynamic Parallel
3 Execution Environment™ (DPEE).

4 117. At least through that correspondence, and no later than February 2013, Mr.
5 Sandstrom provided Microsoft's cloud computing team a copy of the DPEE system description
6 that corresponds with the preferred embodiments disclosed in ThroughPuter's patents. *See* Exhibit
7 20.

8 118. From 2013 through 2016, ThroughPuter's President Mr. Sandstrom continued to
9 correspond with Microsoft's cloud computing team concerning a potential collaboration between
10 Microsoft and ThroughPuter. *See, e.g.,* Exhibit 32.

11 119. For example, in correspondence dated May 12, 2015, in response to a
12 communication from ThroughPuter advising Microsoft its patent portfolio, which at that time
13 included 25 granted US and UK patents, Dr. Burger stated that he had "forwarded
14 [ThroughPuter's] note onto the relevant people looking at any IP and they will contact you directly
15 if interested. Doug." Exhibit 30.

16 120. Shortly thereafter Microsoft filed a patent application directed to the same subject
17 matter as ThroughPuter's technology naming Dr. Burger, to whom Mr. Sandstrom disclosed
18 ThroughPuter's technology, as an inventor. Following the filing of that application, the leader of
19 Microsoft's cloud computing team declined ThroughPuter's offers for collaboration, saying in an
20 August 3, 2015 email to Mr. Sandstrom: "I don't think we can proceed ... we already have our
21 plates very full. Thanks and best wishes." Exhibit 32.

22 121. Thereafter, Mr. Sandstrom reminded Microsoft several times that the solution he
23 shared was the subject of several U.S. and foreign patents.

1 122. In 2015, Mr. Sandstrom had a discussion with Microsoft representative, Dr. Derek
2 Chiou, during which conversation, Dr. Chiou made clear that despite Microsoft's knowledge of
3 ThroughPuter's technology and its patent protection, Microsoft would not be willing to explore
4 any business relationship with ThroughPuter, including with respect to ThroughPuter's intellectual
5 property rights.

6 123. On May 16, 2016, Mr. Sandstrom informed Arun Justus at Microsoft that
7 "ThroughPuter owns the IP rights for the key techniques that will be necessary in realizing any
8 scalable solution for this must-solve challenge." Exhibit 33.

9 124. In addition, on May 19, 2016, Mr. Sandstrom wrote Aki Siponen at Microsoft:
10 "[All] cloud service providers will be facing the fundamental scalability solution . . . and
11 ThroughPuter holds the key patented techniques that will be needed in delivering an effective
12 solution to this must-solve challenge." *Id.*

13 125. On May 17, 2017, via a LinkedIn message, Mr. Sandstrom provided notice to Dr.
14 Burger of the issuance of two ThroughPuter patents. In response to Mr. Sandstrom's message, Dr.
15 Burger, stated "Thanks Mark, appreciate the nice message and will definitely review."

16 126. In a 2017 Microsoft Research publication titled "The Feniks FPGA Operating
17 System for Cloud Computing," Microsoft describes its experimental development of FPGA
18 processors with a hardware based operating system in highly similar terms as ThroughPuter in the
19 above-referenced, earlier patent disclosures as well as publications such as the one presented by
20 Mr. Sandstrom at FPGAworld in 2014. *Compare e.g., Exhibit 43 with Exhibit 25.*

21 127. In that 2017 publication, Dr. Chiou is credited with the "initial exploration which
22 provides valuable experience for the system design." (Exhibit 25 at p. 7, Acknowledgement). Dr.
23
24

1 Chiou is the same individual with whom ThroughPuter spoke in 2015 after Microsoft had
2 acknowledged ThroughPuter's technology and its patent protection.

3 128. On January 2, 2018, Mr. Sandstrom sent a letter (the "January 2018 letter") to Kevin
4 Scott, Chief Technology Officer of Microsoft Corporation, informing Microsoft of the existence
5 of ThroughPuter's patent portfolio and invited a business discussion.

6 129. On September 1, 2018, outside patent counsel for ThroughPuter sent Microsoft
7 Chairman John W. Thompson a letter ("the September 2018 letter") offering Microsoft the
8 opportunity to license or acquire "ThroughPuter's Dynamic Parallel Execution (DPE) technology,
9 protected by a portfolio of forty-eight (48) patents issued and pending worldwide". The
10 September 2018 letter further explained that to the extent Microsoft's current or planned products
11 were built on certain cloud and enterprise computing technologies, "Microsoft needs to obtain a
12 license from ThroughPuter for continued as well as planned future usage of [ThroughPuter's]
13 patents".

14 130. On September 1, 2018, ThroughPuter's outside patent counsel sent the same letter
15 referred to in the preceding paragraph to Microsoft's then-General Counsel.

16 131. Microsoft did not respond to any of the letters referred to in the four preceding
17 paragraphs.

18 132. On information and belief, Microsoft cut off communication with ThroughPuter in
19 or around May 2017 because Microsoft realized that ThroughPuter had been issued patents with
20 earlier priority dates on the technology approach underlying the Microsoft Azure PaaS and its
21 development plans.

22 133. On information and belief, Dr. Burger and Microsoft's other personnel were
23 instructed to stop communicating with ThroughPuter because Microsoft recognized the similarities
24

1 between ThroughPuter's patent protected technology and the execution layer of the Azure PaaS
2 and its development plans.

3 134. Rather than respect ThroughPuter's intellectual property rights, Microsoft made,
4 upon information and belief, the deliberate decision to infringe ThroughPuter's patents. This
5 conscious decision prevents ThroughPuter from being able to compete in the cloud computing
6 space all while Microsoft used the enabling ThroughPuter technology to scale up the world's
7 largest cloud-computing platform. ThroughPuter's patented technology allows Microsoft to claim
8 up to 200-fold technical performance gains.

9 135. Microsoft recognized the novelty of ThroughPuter's technology by attempting to
10 claim ThroughPuter's technology as its own. Specifically, as discussed herein, Microsoft sought
11 and obtained patent protection on the fundamental technologies invented by ThroughPuter years
12 earlier.

13 136. Upon information and belief, to date, Microsoft has not disclosed any information
14 concerning ThroughPuter or its technologies to the USPTO in connection with any Microsoft
15 patent application related to Azure.

16 137. On information and belief, Microsoft has been monitoring ThroughPuter's patent
17 portfolio since at least 2017 and has had actual notice of all ThroughPuter patents that issued since
18 that time.

19 138. In a recently filed petition for *inter partes* review directed to the '242 patent,
20 Microsoft admitted that it has been monitoring ThroughPuter's patent family and has knowledge
21 of the claims asserted herein, which claim priority through the patent application that matured into
22 the '242 patent. Such monitoring provided Microsoft with actual notice of the two patents asserted
23 herein.

3 140. Despite this knowledge, Microsoft has made the deliberate choice to refuse any
4 business arrangement or resolution with ThroughPuter, opting instead to infringe ThroughPuter's
5 patents, thereby preventing ThroughPuter from competing in the market ThroughPuter's
6 technology has enabled.

141. On information and belief, Microsoft's infringement has been continuous,
deliberate and in willful disregard of ThroughPuter's patent rights.

10 (Infringement of U.S. Patent No. 11,150,948)

11 142. ThroughPuter repeats and realleges each and every allegation contained above as
12 though fully set forth herein.

13 143. On October 19, 2021, the United States Patent and Trademark Office duly and
14 legally issued the '948 patent, entitled "Managing Programmable Logic-Based Processing Unit
15 Allocation on a Parallel Data Processing Platform." A copy of the '948 patent is attached as
16 Exhibit 8.

17 144. Mark Sandstrom is the sole and true inventor of the '948 patent.

18 145. ThroughPuter, Inc. owns all right, title, and interest to and in the '948 patent.

19 146. Microsoft infringes claims 1-15 of the '948 patent.

147. As demonstrated in the side-by-side comparison below, claim 1 of ThroughPuter’s ‘948 patent closely matches the claims of Microsoft’s ‘657 patent. Ex. 42. In the chart below, the column on the left shows claim 1 of the ‘948 patent, which claims priority to applications dating back to 2012. The column on the right shows independent claim 8 of Microsoft’s ‘657 patent,

which claims priority to an application filed in 2015. As can be appreciated from this side-by-side comparison, Microsoft obtained a patent on substantially the same technology taught by ThroughPuter's patent applications. However, ThroughPuter's '948 patent is entitled to a priority date that is over two years earlier than Microsoft's '657 patent.

ThroughPuter's U.S. Patent No. 11,150,948 Independent Claim 1	Microsoft's U.S. Patent No. 10,819,657, Independent Claim 8
A method on a parallel data processing platform comprising a plurality of programmable logic-based processing units for accelerating a service comprising at least a first program and a second program, different than the first program, the method comprising:	A method in a data center comprising a pool of acceleration components for accelerating the service comprising at least a first service functionality and a second service functionality, different from the first functionality, the method comprising:
forming a first set of inter-task communication paths connecting a first set of programmable logic-based processing units of the plurality of programmable logic-based processing units into a first multi-stage program instance, wherein the first multi-stage program instance is configured to accelerate the first program, and wherein the first set of programmable logic-based processing units are interconnected using reconfigurable cross-connects;	forming a first graph comprising a first set of interconnected acceleration components, wherein the first graph is configured to accelerate the first service functionality corresponding to the service, and wherein each of the first set of acceleration components comprises a programmable logic array comprising hardware logic blocks interconnected using reconfigurable interconnects;
forming a second set of inter-task communication paths connecting a second set of programmable logic-based processing units of the plurality of programmable logic-based processing units to a second multi-stage program instance, wherein the second multi-stage program instance is configured to accelerate the second program, and wherein the second set of programmable logic-based processing units are interconnected using reconfigurable cross-connects;	forming a second graph comprising a second set of interconnected acceleration components, wherein the second graph is configured to accelerate the second service functionality corresponding to the service, and wherein each of the second set of acceleration components comprises a programmable logic array comprising hardware logic blocks interconnected using reconfigurable interconnects;
in response to a first increased demand for the first multi-stage program instance determined by monitoring a first characteristic corresponding to the first set of programmable logic-based processing units, based on an adaptive optimized resource allocation policy, adding at least one more programmable logic-based processing unit to the first set of programmable logic-based processing units from the plurality of programmable logic-based processing units; and	in response to a first increased demand for hardware acceleration from the service determined by monitoring a first characteristic corresponding to the first set of interconnected acceleration components, based on an allocation policy, a service management component adding at least one more acceleration component to the first set of interconnected acceleration components from the pool of acceleration components;

in response to a second increased demand for the second multi-stage program instance determined by monitoring a second characteristic different from the first characteristic, corresponding to the second set of programmable logic-based processing units, based on the adaptive optimized resource allocation policy, adding at least one more programmable logic-based processing unit to the second set of programmable logic-based processing units from the plurality of programmable logic-based processing units.	and in response to a second increased demand for hardware acceleration from the service determined by monitoring a second characteristic, different from the first characteristic, corresponding to the second set of interconnected acceleration components, based on the allocation policy, the service management component adding at least one more acceleration component to the second set of interconnected acceleration components from the pool of acceleration components.
based on at least the determination, the given target programmable logic core redirects the request to another programmable logic core implementing the designated task.	based on at least the determination, the target hardware acceleration device redirects the request to another hardware acceleration device implementing the designated hardware accelerated service.
configuring the selected field-programmable gate arrays to process a respective processing stage of a respective requesting instance, and	configuring the selected field-programmable gate arrays of the plurality to perform the individual stages of the different functions; and
configuring certain selected field-programmable gate arrays to support communicating, by the task executing on the respective field-programmable gate array, final results to a requesting client over a network in the data processing system.	configuring certain selected field-programmable arrays to communicate final results of the different functions to the requesting server unit components, the certain selected field-programmable arrays communicating the final results to the requesting server unit components over a network in the data processing system.

148. On information and belief, Microsoft believed at the time of filing and still believes that the subject matter claimed in the '657 patent is patent eligible under 35 U.S.C. § 101.

149. On information and belief, Microsoft believed at the time of filing and still believes that the subject matter claimed in the '657 patent is novel under 35 U.S.C. § 102 in view of the references disclosed to and considered by the United States Patent & Trademark Office ("USPTO") during examination of the underlying application.

150. On information and belief, Microsoft believed at the time of filing and still believes that that the subject matter claimed in the '657 patent is non-obvious under 35 U.S.C. § 103 in view of the references Microsoft disclosed to and considered by the USPTO during examination

1 of the underlying application.

2 151. Claim 1 of the '948 patent is representative of the claims infringed by Microsoft
3 and recites:

4 1. A method on a parallel data processing platform comprising a plurality of programmable
5 logic-based processing units for accelerating a service comprising at least a first
6 program and a second program, different than the first program, the method
comprising:

7 forming a first set of inter-task communication paths connecting a first set of programmable
8 logic-based processing units of the plurality of programmable logic-based
9 processing units into a first multi-stage program instance, wherein the first multi-
stage program instance is configured to accelerate the first program, and wherein
the first set of programmable logic-based processing units are interconnected using
reconfigurable cross-connects;

10 forming a second set of inter-task communication paths connecting a second set of
11 programmable logic-based processing units of the plurality of programmable logic-
based processing units to a second multi-stage program instance, wherein the
12 second multi-stage program instance is configured to accelerate the second
program, and wherein the second set of programmable logic-based processing units
are interconnected using reconfigurable cross-connects;

13 in response to a first increased demand for the first multi-stage program instance
14 determined by monitoring a first characteristic corresponding to the first set of
programmable logic-based processing units, based on an adaptive optimized
15 resource allocation policy, adding at least one more programmable logic-based
processing unit to the first set of programmable logic-based processing units from
16 the plurality of programmable logic-based processing units; and

17 in response to a second increased demand for the second multi-stage program instance
18 determined by monitoring a second characteristic different from the first
characteristic, corresponding to the second set of programmable logic-based
19 processing units, based on the adaptive optimized resource allocation policy,
adding at least one more programmable logic-based processing unit to the second
20 set of programmable logic-based processing units from the plurality of
programmable logic-based processing units.

21 152. On information and belief, Azure is implemented in a manner that meets each and
22 every limitation of claim 1 of the '948 patent.

23 153. On information and belief, Azure practices the invention recited in claim 8 of
24 Microsoft's '657 patent.

1 ***1. A method on a parallel data processing platform comprising a plurality of***
 2 ***programmable logic-based processing units for accelerating a service comprising at least***
 3 ***a first program and a second program, different than the first program, the method***
 4 ***comprising:***

5 154. Azure comprises a parallel data processing system with a plurality of
 6 programmable logic-based processing units (FPGA processors) for accelerating multiple
 7 applications such as DNNs, RNNs, compression, web search ranking, speech recognition, and
 8 SDN. In some applications a single application or service (e.g., speech recognition) may have
 9 multiple implementations (e.g., English and German) associated with it. In Azure, applications or
 10 services are often broken down into smaller functional units that are executed (accelerated)
 11 separately.

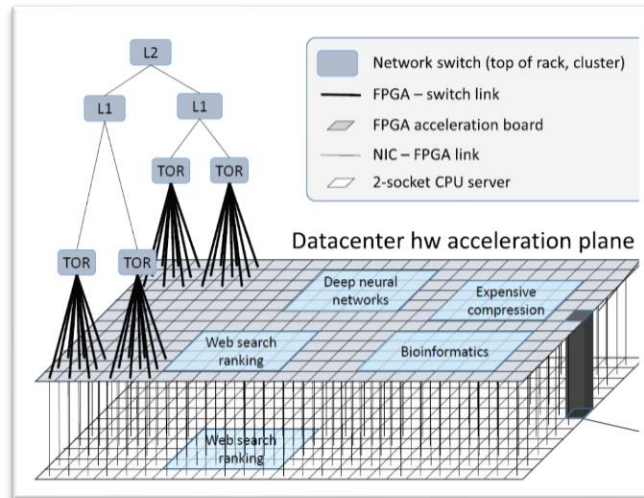
12 ***forming a first set of inter-task communication paths connecting a first set of***
 13 ***programmable logic-based processing units of the plurality of programmable logic-***
 14 ***based processing units into a first multi-stage program instance, wherein the first multi-***
 15 ***stage program instance is configured to accelerate the first program, and wherein the***
 16 ***first set of programmable logic-based processing units are interconnected using***
 17 ***reconfigurable cross-connects;***

18 155. Microsoft uses Azure to form a first set of inter-task communication paths
 19 connecting a first set of programmable logic-based processing units of the plurality of
 20 programmable logic-based processing units into a first multi-stage program instance.

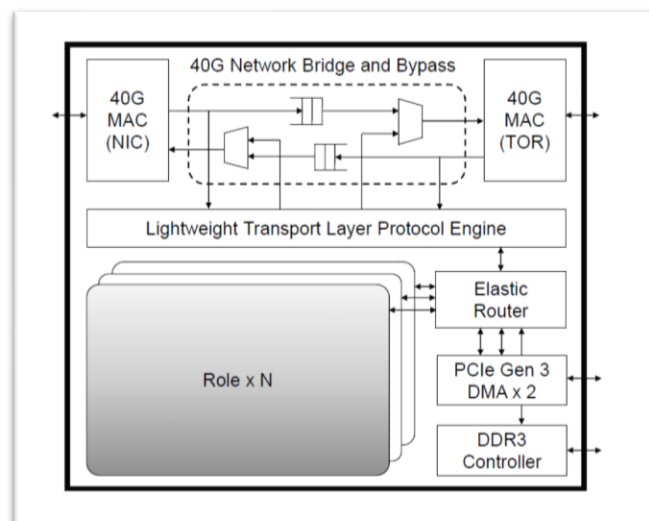
21 156. Azure includes multiple FPGAs, which are programmable logic based processing
 22 units. Microsoft connects a first set of FPGAs into a multi-stage program instance comprised of
 23 multiple FPGA cores. Every FPGA in an Azure datacenter is able to communicate with every
 24 other FPGA in a datacenter through a high-speed, low latency, multiplexer switched fabric. Exs.
 10, 14, 36 *passim*.

157. Each Azure rack in a datacenter is connected by a two-tier switching network L1/L2
 as shown below. Exs. 10, 14, *passim*. The packet communications within each rack and to the top-

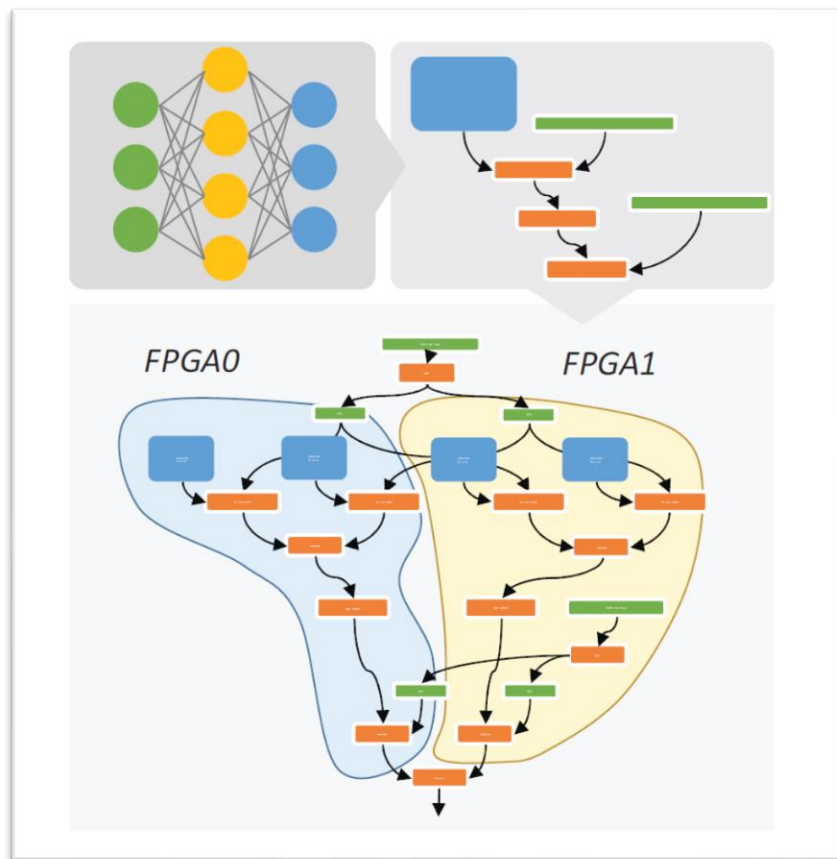
of-rack switch (TOR) are controlled by the FPGA via the high-speed, low-latency multiplexer controlled connections that serve as the Catapult II Shell Components. *Id.*



158. The hardware through which the FPGAs communicate constitutes a first set of inter-task communication paths connecting a first set of programmable logic-based processing units. The figure below shows the “shell” of the Catapult II FPGA. Exs. 10, 14. All of the depicted functionality is performed on the FPGA itself, which is essentially an array of multiplexers and wires that provide high-speed, low latency interconnectivity between logical units. *Id.*, see also Ex. 46 *passim*, Ex. 16 at 5.



159. The FPGA cores for a given program flow can be spread among multiple FPGA stages. Ex. 18 at 10-12, Ex. 37 at 21-22. The group of FPGA cores that are performing the same task constitute a stage. *Id.* In the illustration below (an excerpt of Ex. 18 Fig. 3), the tasks by the DNN model are separated into three types of tasks, illustrated in blue, orange and green. *Id.* Each of tasks is deployed on a FPGA processing core instantiated on one of two FPGA processing stages, labelled FPGA0 and FPGA1. *Id.* The blue cores operate collectively to perform a first task in a first stage and pass the output to another a stage comprised of orange cores. *Id.* The green cores operate collectively to perform a second task in a second stage and pass the output to one of a set of stages comprised of orange cores. *Id.* The orange cores perform a third function and serve as the fourth, fifth and sixth stages, although those stages could be considered a single stage for certain purposes and in certain implementations. *Id.*

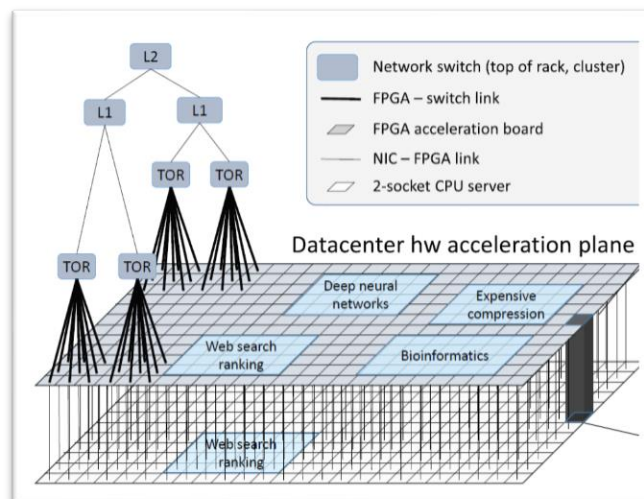


forming a second set of inter-task communication paths connecting a second set of programmable logic-based processing units of the plurality of programmable logic-based processing units to a second multi-stage program instance, wherein the second multi-stage program instance is configured to accelerate the second program, and wherein the second set of programmable logic-based processing units are interconnected using reconfigurable cross-connects;

160. Microsoft uses Azure to form a second set of inter-task communication paths connecting a second set of programmable logic-based processing units of the plurality of programmable logic-based processing units into a first multi-stage program instance.

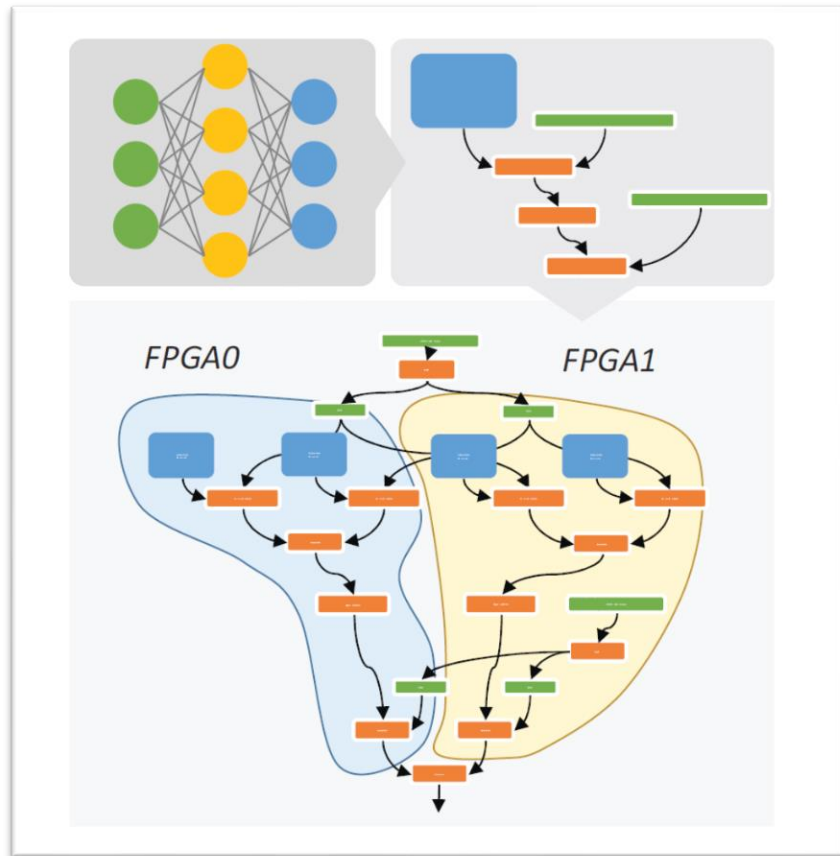
161. Azure includes multiple FPGA processors, which are programmable logic based processing units. Microsoft connects a second set of FPGA processors to form a multi-stage program instance comprised of multiple FPGA cores. Every FPGA in an Azure datacenter is able to communicate with every other FPGA in a datacenter through a high-speed, low latency, multiplexer switched fabric. Exs. 10, 14, 36 *passim*.

162. Each Azure rack in a datacenter is connected by a two-tier switching network L1/L2 as shown below. Exs. 10, 14, *passim*. The packet communications within each rack and to the top-of-rack switch (TOR) are controlled by the FPGA via the high-speed, low-latency multiplexer controlled connections that serve as the Catapult II Shell Components. *Id*.



1 163. The hardware through which the FPGAs communicate constitutes a first set of
2 inter-task communication paths connecting a first set of programmable logic-based processing
3 units.

4 164. The FPGA cores for a given program flow can be spread among multiple FPGA
5 stages. Ex. 18 at 10-12, Ex. 37 at 21-22. The group of FPGA cores that are performing the same
6 task constitute a stage. *Id.* In the illustration below (an excerpt of Ex. 18 Fig. 3), the tasks by the
7 DNN model are separated into three types of tasks, illustrated in blue, orange and green. *Id.* Each
8 of tasks is deployed on a FPGA processing core instantiated on one of two FPGA processing
9 stages, labelled FPGA0 and FPGA1. *Id.* The blue cores operate collectively to perform a first task
10 in a first stage and pass the output to another a stage comprised of orange cores. *Id.* The green
11 cores operate collectively to perform a second task in a second stage and pass the output to one of
12 a set of stages comprised of orange cores. *Id.* The orange cores perform a third function and serve
13 as the fourth, fifth and sixth stages, although those stages could be considered a single stage for
14 certain purposes and in certain implementations. *Id.*



165. Azure simultaneously accelerates multiple multi-stage services or applications. At any given point in time, an Azure datacenter is accelerating multiple DNNs, RNNs, web search ranking, speech recognition, compression, computer vision, speech translation, and other programs. Many of these program or service instances are multi-stage.

in response to a first increased demand for the first multi-stage program instance determined by monitoring a first characteristic corresponding to the first set of programmable logic-based processing units, based on an adaptive optimized resource allocation policy, adding at least one more programmable logic-based processing unit to the first set of programmable logic-based processing units from the plurality of programmable logic-based processing units; and

166. In response to a first increased demand for the first multi-stage program instance determined by monitoring a first characteristic corresponding to the first set of programmable logic-based processing units, based on an adaptive optimized resource allocation policy, Azure

1 adds at least one more programmable logic-based processing unit to the first set of programmable
2 logic-based processing units from the plurality of programmable logic-based processing units.

3 167. As Microsoft explained in a white paper discussing the Catapult II architecture,
4 “[a]s demand for a service grows or shrinks, a global manager grows or shrinks the pools
5 correspondingly.” Ex. 10 at 2, see also Ex. 14 at 4. The pools of cores among and between
6 different FPGAs assigned to each multi-stage program instance is controlled by the Catapult II
7 Shell Components under the supervision and control of the Service Manager and Resource
8 Manager. Ex. 10 at 10. As discussed above, load data, which is a first characteristic as claimed, is
9 monitored and tracked at each FPGA core group. See Ex.41. That data is published to other FPGA
10 core groups in the datacenter. Based on this monitoring, Azure expands and contracts the number
11 of cores assigned to each group over time.

12 *in response to a second increased demand for the second multi-stage program instance*
13 *determined by monitoring a second characteristic different from the first characteristic,*
14 *corresponding to the second set of programmable logic-based processing units, based on*
15 *the adaptive optimized resource allocation policy, adding at least one more*
16 *programmable logic-based processing unit to the second set of programmable logic-*
17 *based processing units from the plurality of programmable logic-based processing units.*

18 168. In response to a second increased demand for the second multi-stage program
19 instance determined by monitoring a second characteristic different from the first characteristic,
20 corresponding to the second set of programmable logic-based processing units, based on the
21 adaptive optimized resource allocation policy, Azure adds at least one more programmable logic-
22 based processing unit to the second set of programmable logic-based processing units from the
23 plurality of programmable logic-based processing units.

24 169. Multiple instances of such multi-stage programs are executed in similar manner
described above in connection with the first multi-stage program instance. That discussion is
incorporated herein by reference.

1 170. As described in Microsoft’s own patent filings, “the load data [need] not [be] stored
2 as a single discrete value. For example, the load data 414 may include a plurality of data types and
3 values. Such as, for example, a number of queued requests for the target hardware acceleration
4 device, a recent processing time of a previous request, an estimate based on the number and type
5 of requests in the queue, a total size of the queued requests, a round trip time for receiving
6 responses to requests, and any other suitable types of data that may indicate a load of the target
7 hardware acceleration device.” Ex.41, 472 Patent at 8:22-24. Each of the foregoing represents a
8 second characteristic different than the first characteristic as claimed.

9 171. On information and belief, Azure uses a plurality of these different data types to
10 reallocate cores among the multi-stage program instances in a given data center on any given day.
11 Any two such multi-stage program instances being rebalanced or reallocated according to given
12 load data types may constitute the first and second multi-stage program instances.

13 172. Upon information and belief, in accordance with 35 U.S.C. § 287, Microsoft has
14 had actual notice and knowledge of the ’948 patent no later than its issuance.

15 173. Microsoft continues without license to make, use, import, offer for sale, and/or sell
16 in the United States services or products that infringe the ’948 patent including specifically
17 Microsoft Azure and its cloud computing functionalities.

18 174. Microsoft has directly and indirectly infringed and continues to directly and
19 indirectly infringe the ’948 patent by engaging in acts constituting infringement under 35 U.S.C.
20 § 271(a) including but not necessarily limited to one or more of making, using, selling and offering
21 to sell, in this District and elsewhere in the United States, and importing into the United States, the
22 Microsoft Azure platform or components and services thereof.

23 175. Microsoft’s continued infringement of the ’948 patent is knowing, intentional, and
24

1 willful.

2 176. Microsoft has had knowledge of and notice of the chain of applications underlying
3 the '948 patent since at least May 2015, when Microsoft representatives acknowledged receipt of
4 an email from ThroughPuter disclosing several of ThroughPuter's U.S. patent applications and in
5 any case no later than January 2018 when ThroughPuter brought ThroughPuter's patent portfolio
6 to the attention of Microsoft in the January 2018 letter, and despite this knowledge continues to
7 commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and
8 above, this infringement has been willful.

9 177. Microsoft has been actively monitoring ThroughPuter's patent portfolio since at
10 least the time it was served with the complaint in the Virginia Action and has gained actual
11 knowledge of the patents asserted herein.

12 178. Microsoft actively, knowingly, and intentionally has induced, or has threatened to
13 induce, infringement of the '948 patent through a range of activities.

14 179. First, on information and belief, Microsoft has induced infringement by, with
15 knowledge of the '948 patent, controlling the design and development of, offering for sale, and
16 selling the services of the Azure platform with the knowledge and specific intent that its customers
17 will use the Azure platform to infringe the '948 patent by executing the system operations and
18 utilizing the system components to perform dynamic resource management of the pool of Azure
19 processing resources on behalf of application programs of customers of the Azure platform cloud
20 computing services.

21 180. Second, on information and belief, Microsoft has, with knowledge of the '948
22 patent, induced infringement by its customers through the dissemination of promotional,
23 marketing, and tutorial materials relating to the Azure platform with the knowledge and specific
24

1 intent that its customers will use the Azure platform to infringe the '948 patent by executing the
2 system operations and utilizing the system components to perform dynamic resource management
3 of the pool of Azure processing resources on behalf of application programs of customers of the
4 Azure platform cloud computing services.

5 181. Third, on information and belief, Microsoft has, with knowledge of the '948 patent,
6 induced infringement by its customers through the creation and online posting of tutorial and
7 "how-to" materials for the Azure platform and/or its individual components in the United States
8 with the knowledge and specific intent that its customers will use the Azure platform to infringe
9 the '948 patent by executing the system operations and utilizing the system components to perform
10 dynamic resource management of the pool of Azure processing resources on behalf of application
11 programs of customers of the Azure platform cloud computing services.

12 182. Fourth, on information and belief, Microsoft has, with knowledge of the '948
13 patent, induced infringement through the distribution of other instructional materials, product
14 manuals, and technical materials with the knowledge and the specific intent to encourage and
15 facilitate its customers' infringing use of the Azure platform.

16 183. Microsoft has engaged in the above activities with knowledge of the '948 patent
17 and with the specific intent to encourage and cause infringement by its customers, as shown by the
18 allegations set forth above.

19 184. Microsoft has contributed to, or has threatened to contribute to, the infringement by
20 its customers of the '948 patent by, without authority, selling and offering to sell within the United
21 States cloud computing services and customer support services for practicing the claimed
22 invention of the '948 patent, including at least the Azure platform as a whole and/or the individual
23 components of the Azure platform. When, for example, the Azure platform is used by Microsoft's
24

1 customers for the various cloud computing services Microsoft offers, the Azure system operations
2 and system components are used to perform the claimed dynamic resource management of the
3 pool of Azure processing resources on behalf of application programs of customers of the Azure
4 cloud platform, thereby infringing the '948 patent. In such settings, the Azure platform and/or its
5 individual components, supplied by Microsoft, constitute at least a material part of the claimed
6 invention of the '948 patent. Microsoft's infringement of the '948 patent has injured ThroughPuter
7 in its business and property rights.

8 185. For example, Microsoft supplies certain customers with an Azure Stack Edge Pro
9 FPGA, which provides customers with a Hardware-as-a-service solution. *See*
10 <https://docs.microsoft.com/en-us/azure/databox-online/azure-stack-edge-overview>. The Azure
11 Stack Edge Pro FPGA can be used for "rapid Machine Learning (ML) inferencing at the edge and
12 preprocessing data before sending it to Azure."

13 186. The Azure Stack Edge Pro FPGA is a component of a patented system that is used
14 in practicing a patent process that constitutes at least a material part of the invention. The Azure
15 Stack Edge Pro FPGA had no substantial non-infringing use at least because any non-infringing
16 use for machine learning would be *de minimis*.

17 187. Microsoft's infringement of the '948 patent has been and is deliberate and willful
18 and constitutes egregious misconduct. On information and belief, despite actual knowledge of the
19 '948 patent and numerous related patents and applications since at least May 2015, Microsoft
20 continued to develop and offer its infringing products and services. In developing and offering its
21 products and services, Microsoft has been willfully blind to this ongoing infringement.

22 188. Pursuant to 35 U.S.C. § 284, ThroughPuter is entitled to recover monetary damages
23 for the injuries arising from Microsoft's willful infringement in an amount to be determined at
24

trial. Microsoft's infringement of the '948 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

189. Microsoft's infringement of the '948 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 2

(Infringement of U.S. Patent No. 11,036,556)

190. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

191. On June 15, 2021, the United States Patent and Trademark Office duly and legally issued the '556 patent, entitled "Concurrent Program Execution Optimization." A copy of the '556 patent is attached as Exhibit 1.

192. Mark Sandstrom is the sole and true inventor of the '556 patent.

193. ThroughPuter, Inc. owns all right, title and interest to and in the '556 patent.

194. Microsoft infringes claims 1-8 of the '556 patent.

195. As demonstrated in side-by-side comparison, claim 1 of ThroughPuter's '556 patent closely matches claim 17 of Microsoft's '392 patent. Ex. 47. The column on the left shows claim 1 of the '556 patent, which claims priority to applications dating back to 2012. The column on the right shows independent claim 17 of Microsoft's '392 patent, which claims priority to an application filed in 2015. As can be appreciated from this side-by-side comparison, Microsoft obtained a patent on substantially the same technology taught by ThroughPuter's patent

1 applications. However, ThroughPuter's '556 patent is entitled to a priority date that is more than
2 two years earlier than Microsoft's.

ThroughPuter's U.S. Patent No. 11,556,556, Independent Claim 1	Microsoft's U.S. Patent No. 10,296,392, Independent Claim 17
A method performed in a data processing system, the method comprising:	A method performed in a data processing system, the method comprising:
receiving, by hardware logic and/or software logic, requests to perform different tasks on behalf of instances of a plurality of programs managed by a data processing system;	receiving requests to perform different functions on behalf of tenant functionality that executes on requesting server unit components of the data processing system;
identifying, by the hardware logic and/or software logic for each of the instances, communication interdependencies between different processing stages of a set of processing stages of the respective instance;	parsing the different functions to identify programmatic calls between different stages of the different functions;
based on conditions in the data processing system, dynamically varying, by the hardware logic and/or software logic, structures of field-programmable gate arrays used to process different tasks of the instances of the plurality of programs, the structures being dynamically varied by	based at least on conditions in the data processing system, dynamically varying structures of field-programmable gate arrays used to implement different invocations of the different functions, the structures being dynamically varied by:
identifying available field-programmable gate arrays of the data processing system that are available to process different processing stages of requesting instances of respective programs,	identifying available field-programmable arrays of the data processing system that are available to implement the different stages of the different functions;
based at least on the conditions in the data processing system, identifying selected field-programmable gate arrays from the available field-programmable gate arrays to execute the different processing stages of the requesting instances of the respective programs,	based at least on the conditions in the data processing system, identifying selected field-programmable gate arrays from the available field-programmable gate arrays to perform individual stages of the different functions;
configuring the selected field-programmable gate arrays to process a respective processing stage of a respective requesting instance, and	configuring the selected field-programmable gate arrays of the plurality to perform the individual stages of the different functions; and
configuring certain selected field-programmable gate arrays to support communicating, by the task executing on the respective field-programmable gate array, final results to a requesting client over a network in the data processing system.	configuring certain selected field-programmable arrays to communicate final results of the different functions to the requesting server unit components, the certain selected field-programmable arrays communicating the final results to the requesting server unit components over a network in the data processing system.

1 196. On information and belief, Microsoft believed at the time of filing and still believes
2 that the subject matter claimed in the '392 patent is patent eligible under 35 U.S.C. § 101.

3 197. On information and belief, Microsoft believed at the time of filing and still believes
4 that the subject matter claimed in the '392 patent is novel under 35 U.S.C. § 102 in view of the
5 references disclosed to and considered by the United States Patent & Trademark Office
6 ("USPTO") during examination of the underlying application.

7 198. On information and belief, Microsoft believed at the time of filing and still believes
8 that that the subject matter claimed in the '392 patent is non-obvious under 35 U.S.C. § 103 in
9 view of the references Microsoft disclosed to and considered by the USPTO during examination
10 of the underlying application.

11 199. Claim 1 of the '556 patent is representative of the claims infringed by Microsoft
12 and recites:

13 1. A method performed in a data processing system, the method comprising:

14 receiving, by hardware logic and/or software logic, requests to perform different
15 tasks on behalf of instances of a plurality of programs managed by a data
16 processing system;

17 identifying, by the hardware logic and/or software logic for each of the instances,
18 communication interdependencies between different processing stages of a
19 set of processing stages of the respective instance;

20 based on conditions in the data processing system, dynamically varying, by the
21 hardware logic and/or software logic, structures of field-programmable gate
22 arrays used to process different tasks of the instances of the plurality of
23 programs, the structures being dynamically varied by

24 identifying available field-programmable gate arrays of the data processing system
25 that are available to process different processing stages of requesting
26 instances of respective programs,

27 based at least on the conditions in the data processing system, identifying selected
28 field-programmable gate arrays from the available field-programmable gate
29 arrays to execute the different processing stages of the requesting instances
30 of the respective programs,

1 configuring the selected field-programmable gate arrays to process a respective
2 processing stage of a respective requesting instance, and

3 configuring certain selected field-programmable gate arrays to support
4 communicating, by the task executing on the respective field-programmable
gate array, final results to a requesting client over a network in the data
processing system.

5 200. On information and belief, Azure is implemented in a manner that meets each and
6 every limitation of claim 1 of the '556 patent.

7 201. On information and belief, Azure practices the invention recited in claim 17 of
8 Microsoft's '392 patent (Ex. 47).

9 ***1. A method performed in a data processing system, the method comprising:***

10 202. Azure is a data processing system.

11 ***receiving, by hardware logic and/or software logic, requests to perform different tasks
on behalf of instances of a plurality of programs managed by a data processing system;***

12 203. Azure receives by hardware logic/and or software logic, requests to perform
13 different tasks on behalf of instances of a plurality of programs managed by a data processing
14 system.

15 204. Azure, e.g. via Catapult II Shell Components instantiated on FPGAs, receives
16 requests to accelerate services from at least the CPUs by hardware logic and/or software logic that
17 are positioned "behind" the FPGAs in the Catapult II network. Ex. 15 at 54. On information and
18 belief, in Azure's current form a hardware microservice running on an FPGA may also initiate a
19 request for acceleration.

20 ***identifying, by the hardware logic and/or software logic for each of the instances,
communication interdependencies between different processing stages of a set of
21 processing stages of the respective instance;***

22 205. Azure identifies, by the hardware logic and/or software logic for each of the
23 instances, communication interdependencies between processing stages of a set of processing
24

1 stages of the respective instance.

2 206. In Azure, under the supervision of the RMs the SM components instantiated on the
3 CPUs and/or FPGA shells determine how to break a requested service down into steps or sub-task
4 components. Ex. 35 at 119, Ex. 18 at 11. Those sub-tasks are assigned to various processing cores
5 on the FPGA fabric. As part of this process, the SM components determine what each sub-task
6 will provide as inputs and outputs, which involves identification of communication
7 interdependencies between the stages (e.g. which stage needs to communicate with which other
8 stage).

9 ***based on conditions in the data processing system, dynamically varying, by the hardware***
10 ***logic and/or software logic, structures of field-programmable gate arrays used to process***
different tasks of the instances of the plurality of programs,

11 207. Based on conditions in the data processing system, Azure dynamically varies, by
12 the hardware logic and/or software logic, structures of field-programmable gate arrays used to
13 process different tasks of the instances of the plurality of programs.

14 208. Azure's AC, RM, SMs and FMs works in conjunction to assign tasks of a multi-
15 stage program instance to a group of reconfigurable FPGA processors. When a core is reallocated
16 to perform a different task (for instance speech recognition instead of compression), the FPGA
17 circuit is reconfigured accordingly by loading a corresponding image file or role. Ex. 10 at 7, Ex.
18 35 at 116. Such reconfiguration varies the structure of the logic blocks of the FPGA.

19 ***the structures being dynamically varied by identifying available field-programmable gate***
20 ***arrays of the data processing system that are available to process different processing***
stages of requesting instances of respective programs,

21 209. In Azure the process of dynamically varying structures of available field-
22 programmable gate arrays of the data processing system, includes identifying PFAs that are
23 available to process different processing stages of requesting instances of respective programs.

24 210. As discussed above in connection with previous counts and/or in the section of this

complaint entitled Microsoft's Infringing Cloud Computing Architecture, Azure, e.g. via the Catapult II system, monitors load data of various types and reallocates FPGA processing cores accordingly. This reallocation includes identifying available FPGAs. For instance, unused or available cores are allocated to "join" groups that are experiencing particularly high demand.

based at least on the conditions in the data processing system, identifying selected field-programmable gate arrays from the available field-programmable gate arrays to execute the different processing stages of the requesting instances of the respective programs,

211. Azure, based at least on the conditions in the data processing system, identifies selected field-programmable gate arrays from the available field-programmable gate arrays to execute the different processing stages of the requesting instances of the respective programs.

212. As discussed above, after such a reallocation involving a multi-stage stage program instance (e.g. an instance of a DNN service), the SM(s) in cooperation with the FMs determine which sub-tasks will be executed on which cores of which newly allocated FPGA blocks. These assignments are made based in part on the load data reported concerning processing tasks placed on each FPGA or block thereof, which data comprises the conditions in the data processing system.

configuring the selected field-programmable gate arrays to process a respective processing stage of a respective requesting instance, and

213. Azure configures the selected field-programmable gate arrays to process a respective processing stage of a respective.

214. As discussed above, when cores are allocated for execution of a new function that requires execution of a new function, FPGA core blocks may be re-imaged as needed. That re-imaging may be partial (only some cores on the FPGA) or complete (all processing cores on the FPGA). Typically, the Catapult II Shell Components are not re-imaged during this process. Such reimaging constitutes configuring the FPGA gate arrays as claimed.

1 *configuring certain selected field-programmable gate arrays to support communicating,*
2 *by the task executing on the respective field-programmable gate array, final results to a*
3 *requesting client over a network in the data processing system.*

4 215. Azure configures selected field-programmable gate arrays to support
5 communicating, by the task executing on the respective field-programmable gate array, final
6 results to a requesting client over a network in the data processing system.

7 216. The head component of each group of cores (each instantiated on an FGPA) is
8 configured under the supervision of the RM and the SM components to communicate the final
9 result of the acceleration request to the requesting client, hosted e.g. on the CPU which is local to
10 the FPGA launching the acceleration request. On information and belief, in its current form Azure,
11 e.g. via Catapult II, permits requesting clients to reside on FPGAs (e.g., as hardware accelerated
12 microservices). Azure also allows external CPUs to communicate with FPGAs.

13 217. Upon information and belief, in accordance with 35 U.S.C. § 287, Microsoft has
14 had actual notice and knowledge of the '556 patent no later than its issuance.

15 218. Microsoft has been actively monitoring ThroughPuter's patent portfolio since at
16 least the time it was served with the complaint in the Virginia Action and has gained actual
17 knowledge of the patents asserted herein.

18 219. Microsoft continues, without license, to make use, offer for sale, import and/or sell
19 in the United States services or products that infringe the '556 patent including specifically
20 Microsoft Azure and its cloud computing functionalities.

21 220. Microsoft has directly infringed and continues to directly infringe the '556 patent
22 by engaging in acts constituting patent infringement under 35 U.S.C. § 271(a) including but not
23 necessarily limited to one or more of making, using, selling and offering to sell, in this District and
24 elsewhere in the United States, and importing into the United States, the Microsoft Azure platform
or components and services thereof.

1 221. Microsoft’s continuing infringement of the ’556 patent is knowing, intentional, and
2 willful.

3 222. Microsoft has had knowledge of and notice of the chain of applications underlying
4 the ’556 patent since at least May 2015, when Microsoft representatives acknowledged receipt of
5 an email from ThroughPuter disclosing several of ThroughPuter’s U.S. patent applications and in
6 any case no later than January 2018 when ThroughPuter brought ThroughPuter’s patent portfolio
7 to the attention of Microsoft in the January 2018 letter, and despite this knowledge continues to
8 commit the aforementioned infringing acts – and despite this knowledge continues to commit the
9 aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this
10 infringement has been willful.

11 223. Microsoft actively, knowingly, and intentionally has induced, or has threatened to
12 induce, infringement of the ’556 patent under 35 U.S.C. § 271(b) and (c) through a range of
13 activities.

14 224. First, on information and belief, Microsoft has induced infringement by, with
15 knowledge of the ’556 patent, controlling the design and development of, offering for sale, and
16 selling the services of the Azure platform with the knowledge and specific intent that its customers
17 will use the Azure platform to infringe the ’556 patent by executing the system operations and
18 utilizing the system components to perform dynamic resource management of the pool of Azure
19 processing resources on behalf of application programs of customers of the Azure platform cloud
20 computing services.

21 225. Second, on information and belief, Microsoft has, with knowledge of the ’556
22 patent, induced infringement by its customers through the dissemination of promotional,
23 marketing, and tutorial materials relating to the Azure platform with the knowledge and specific
24

1 intent that its customers will use the Azure platform to infringe the '556 patent by executing the
2 system operations and utilizing the system components to perform dynamic resource management
3 of the pool of Azure processing resources on behalf of application programs of customers of the
4 Azure platform cloud computing services.

5 226. Third, on information and belief, Microsoft, with knowledge of the '556 patent,
6 has induced infringement by its customers through the creation and online posting of tutorial and
7 "how-to" materials for the Azure platform and/or its individual components in the United States
8 with the knowledge and specific intent that its customers will use the Azure platform to infringe
9 the '556 patent by executing the system operations and utilizing the system components to perform
10 dynamic resource management of the pool of Azure processing resources on behalf of application
11 programs of customers of the Azure platform cloud computing services.

12 227. Fourth, on information and belief, Microsoft has, with knowledge of the '556
13 patent, induced infringement through the distribution of other instructional materials, product
14 manuals, and technical materials with the knowledge and the specific intent to encourage and
15 facilitate its customers' infringing use of the Azure platform.

16 228. Microsoft has engaged in the above activities with knowledge of the '556 patent
17 and with the specific intent to encourage and cause infringement by its customers, as shown by the
18 allegations set forth above.

19 229. Microsoft has contributed to, or has threatened to contribute to, the infringement by
20 its customers of the '556 patent by, without authority, selling and offering to sell within the United
21 States cloud computing services and customer support services for practicing the claimed
22 invention of the '556 patent, including at least the Azure platform as a whole and/or the individual
23 components of the Azure platform. When, for example, the Azure platform is used by Microsoft's
24

1 customers for the various cloud computing services Microsoft offers, the Azure system operations
2 and system components are used to perform the claimed dynamic resource management of the
3 pool of Azure processing resources on behalf of application programs of customers of the Azure
4 cloud platform, thereby infringing the '556 patent. In such settings, the Azure platform and/or its
5 individual components, supplied by Microsoft, constitute at least a material part of the claimed
6 invention of the '556 patent.

7 230. Microsoft's infringement of the '556 patent has injured ThroughPuter in its
8 business and property rights.

9 231. Microsoft's infringement of the '556 patent has been and is deliberate and willful
10 and constitutes egregious misconduct. Upon information and belief, despite actual knowledge of
11 the '556 patent and numerous related patents and applications since at least May 2015, Microsoft
12 continued to develop and offer its infringing products and services. In developing and offering its
13 products and services, Microsoft has been willfully blind to this ongoing infringement.

14 232. Pursuant to 35 U.S.C. § 284, ThroughPuter is entitled to recover monetary damages
15 for the injuries arising from Microsoft's willful infringement in an amount to be determined at
16 trial. Microsoft's infringement of the '556 patent has caused irreparable harm to ThroughPuter and
17 will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by
18 this Court.

19 233. Microsoft's infringement of the '556 patent is exceptional and entitles
20 ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. §
21 285.

PRAYER FOR RELIEF

WHEREFORE, ThroughPuter respectfully requests that the Court enter judgment against Microsoft as follows:

A. An adjudication that Microsoft has infringed one or more claims of the '948 and '556 patents;

B. An order permanently enjoining Microsoft from further infringement of the '948 and '556 patents;

C. An award of damages pursuant to 35 U.S.C. § 284;

D. An order that the damages award be increased up to three times the actual amount assessed, pursuant to 35 U.S.C. § 284;

E. An award to ThroughPuter of its costs, pre- and post-judgment interest, and reasonable expenses to the fullest extent permitted by law;

F. A declaration that this case is exceptional pursuant to 35 U.S.C. § 285, and an award of attorneys' fees and costs; and

G. An award to ThroughPuter of such other and further relief as this Court deems just and proper.

DEMAND FOR JURY TRIAL

Pursuant to Rule 38(b) of the Federal Rules of Civil Procedure, ThroughPuter hereby demands a trial by jury on all issues so triable.

1 DATED this 13th day of April, 2022.

2 LOWE GRAHAM JONES^{PLLC}

3 s/Lawrence D. Graham

4 Lawrence D. Graham, WSBA No. 25402

5 Graham@LoweGrahamJones.com

6 1325 Fourth Avenue, Suite 1130

7 Seattle, Washington 98101

8 T: 206.381.3300

9 F: 206.381.3301

10 *Attorneys for Plaintiff ThroughPuter, Inc.*